

Orthogonal Recursive Bisection on the GPU for Accelerated Load Balancing in Large N-Body Simulations

Bachelor Thesis

Andrin Rehmann

18-916-882

Department for Informatics
University of Zürich

Prof. Dr. Michael Böhlen
Supervisor: Douglas Potter
14.08.2022

Abstract

Large simulations with billions of particles are used to understand the universe. Commonly, the fast multipole method is employed to improve the runtime, where a space partitioning binary tree data structure is required. The structure is generated using the Orthogonal Recursive Bisection method and at the same time used as a load balancing strategy to leverage the hardware of super computers optimally.

In this paper, a GPU accelerated version of ORB is proposed and implemented in the CUDA programming language. The implementation manages to maintain a consistent runtime during the entire execution of the ORB method, regardless of the increasing complexity as fragmentation grows with the tree depth and the number of leaf cells. Performance measurements showed a speedup by a factor of 5.4 over a fully parallelized and optimized CPU version.

Contents

1	Introduction	4
2	Background	5
2.1	Fast Multipole Expansion	5
2.2	ORB Compared to other Space Partitioning Methods	6
2.3	Rise of the GPU	6
2.4	Related Work	7
3	Orthogonal Recursive Bisection (ORB)	8
3.1	Target Data and Notation	8
3.2	Memory and Workload Balancing	9
3.3	ORB	10
3.3.1	Algorithm	11
3.4	By Example	12
3.5	Root finding	15
3.5.1	Bisection Method	16
3.5.2	Edge Cases	17
3.5.3	Runtime Analysis	18
3.6	Partition Algorithm	18
4	Theoretical Analysis	20
4.1	General Memory Model	20
4.2	Supercomputers	21
4.3	Roofline Performance Model	22
4.3.1	Estimating Flops	22
4.3.2	Estimating Arithmetic Intensity	23
4.3.3	The Plots	24
4.3.4	Empirical Verification	27
4.4	Runtime Estimates	28
4.4.1	CPU Version	29
4.4.2	GPU Counting	29
4.4.3	GPU Counting and Partitioning	29
4.4.4	Plugin Values	30
4.5	Conclusion	30
5	CPU Implementation	32
5.1	MDL and PKDGrav	32
5.2	Particles	32
5.3	Cell	33
5.3.1	Heap Conditions	33
5.3.2	Class	34
5.4	Mapping Cells to Particles	35
5.5	Services	35
5.5.1	Init and Finalize Service	35

5.5.2	Count Left Service	35
5.5.3	Count	36
5.5.4	Partition	37
5.5.5	Make Axis	38
5.6	Parallel Schedule	38
6	GPU Implementation	41
6.1	Relevant CUDA Concepts	41
6.1.1	Warps	42
6.1.2	Memory	42
6.1.3	Memory Access Patterns	43
6.1.4	Asynchronous Operations	43
6.1.5	Synchronization	43
6.2	Streams	44
6.3	Memory Management	44
6.4	GPU Accelerated Count Left	44
6.4.1	Schedule	44
6.4.2	Service	46
6.4.3	Kernel Code	47
6.5	Improved GPU Accelerated Count Left	49
6.5.1	Schedule	50
6.5.2	Kernel Code	52
6.6	Improved GPU Accelerated Count Left with Warp Level Primitives	53
6.7	GPU Accelerated Partitioning	56
6.7.1	Schedule	57
6.7.2	Partitioning a Single Cell	58
6.7.3	Partitioning Multiple Cells	59
7	Performance Analysis of ORB	60
7.1	Scaling	61
8	Conclusion	62

1 Introduction

The N-Body technique has been used for decades to simulate the Universe to compare theory with observations. It uses *particles* to represent objects of a certain mass and computes the interacting forces by applying the gravity equation. As the forces operate over an infinite distance it is necessary to consider all pairwise interactions, making a naive implementation $\mathcal{O}(N^2)$ where N is the number of particles in the system. This does not scale with large N 's, in fact the problem becomes computationally infeasible at some point.

A common solution is to partition the space, in which the particles are contained, into a set of subspaces using the Orthogonal Recursive Bisection (ORB) algorithm. These subspaces, or cells, are then stored in a space partitioning tree data structure (SPTDS). The data structure is then leveraged in combination with the Fast Multipole Method (FMM) to speed up the particle simulation to $\mathcal{O}(N)$. In order to balance the load of work and memory across nodes and processors in a computing system, a load balancing is necessary. The same tree, which is used for FMM can be leveraged to generate groupings of the particles which are then distributed among processing units. Building the SPTDS with ORB uses a significant percentage of the overall simulation time, [Stadel, 2001] where as of now this was done on the CPU, not leveraging GPU acceleration.

In this paper I established an upper limit for the speedup of a GPU ORB implementation over its fully parallelized CPU only counterpart. In order to do so, I proposed a runtime estimate for both versions. As hardware specific details and some knowledge about compilers and Assembly instructions are crucial for a reliable estimate, I explored the theoretical foundations of memory bandwidth and processing power.

Based on encouraging figures from the runtime estimates I proceeded with the actual implementation. Using the machine dependent layer (MDL) from PKDGRAV [Stadel, 2001], which is used to distribute workload among cores and processors, I implemented a fully parallelized CPU version of ORB. In preparation of advancing parts of the code to a GPU accelerated implementation, I have summarized the most important and relevant concepts of CUDA, a popular graphics programming language developed by NVIDIA. Compared to traditional C++ code, many hardware related limitations have to be considered and a fundamentally different paradigms is required.

With the gained knowledge I implemented the most performance critical part of ORB using CUDA [NVIDIA, 2022e]. The identified part is essentially a map and reduce on the particles coordinates. I iteratively increased the performance of the kernels and describe each version in detail.

Moreover I propose a possible method to accelerate the partitioning part of ORB with specific details. The proposed implementation could possibly further decrease the runtime of ORB when implemented.

Finally, I have performed a runtime analysis on Piz Daint [CSCS, 2022], a powerful super computer in Switzerland. A significant speedup of the GPU accelerated ORB over its CPU counterpart was found.

2 Background

Simulations of the universe mock behavior of the real observed data and allow to gain insight into the process behind certain structures [Stadel, 2001]. Generated datasets with similar distributions to reality are used to seed the simulations. We model a collision less fluid and use "particles" to sample the mass distribution, where each of them represents a mass point. Numerical integration methods then compute the gravitational forces on each object. Originally developed by Joachim Stadel during his PhD at University of Washington [Stadel, 2001], PKDGrav is a fully parallel N-Body code which is now maintained and improved by the ICS at UZH [UZH, 2022]. PkDGrav, as well as most state of the art simulation softwares, use a space partitioning tree data structure (SPTDS), which is used to partition a space into subspaces. Each cell (node), in the tree represents a subspace, and the children of the node represent the subspaces that are contained within the node's subspace. The multipole expansion is then applied to each cell of the SPTDS and the fast multipole method increases the runtime from $O(N^2)$ to $O(N)$ where N is the total number of particles. [Stadel, 2001]

2.1 Fast Multipole Expansion

In mathematics, a multipole is a concept used to approximate the shape of a distribution. In this case the body consists of the particles contained within the subspace, or volume of a cell. The multipole expansion is the mathematical series that results from this approximation. The resulting series can simplify the original body, by choosing a suitable precision.

To keep force integration within a reasonable error margin, in most cases it is not necessary to compute all N to N interactions. The fast multipole method can compute gravity in $O(N)$ runtime by leveraging the SPTDS. [Stadel, 2001]

On a more intuitive level, one can imagine a large cluster of objects far away from planet Earth. To compute the interactions between every single object of the cluster and Earth, the gravity equation needs to be solved $O(N^2)$ times. However, since the object is reasonably far away, it does not make sense to compute every interaction. Instead, the cluster can be summarized into a single object and compute the gravity acting between planet Earth and the group. Consequently, objects within our own solar system would be too close to planet Earth to approximate its effects using a multipole. In such cases, it is necessary to perform all computations, however they can be hidden within a constant and do not affect the total runtime. Note that this example does not directly translate to the simulations, as a particle is not equivalent to a planet or a star. More on that in section 3.1.

The distance between the center of mass of a cell and its most distant corner correlates with the error ratio [Stadel, 2001], which in turn influences how many timesteps have to be performed. Of the family of rectangular cuboids, the cube is the shape where the average distance between any point and its most distant corner point is the smallest.

2.2 ORB Compared to other Space Partitioning Methods

A number of methods exist to create SPTDS, however since the generated tree should be usable for load balancing and FMM, it is necessary to have a load in terms of either work or memory on each computing unit. Naturally the reasonable solution is to use the exact same particles contained within cells for load balancing and FMM, which rules out many space partitioning algorithms such as octrees. In octrees the resulting subspaces are equal of size, but not equal in terms of memory and workload of the contained objects. The k-d tree construction method however is strikingly similar to the ORB method but there is a crucial difference in the construction process: The Orthogonal Recursive Bisection method constructs a SPTDS where all resulting cells represent subspaces, or volumes, similar to cubes. K-d construction methods generate SPTDS as well, but shapes of the resulting subspaces can be strongly skewed and are thus ill suited for FMM. Both algorithms start with a single cell and iteratively construct the SPTDS by searching for axis aligned planes cutting the volume of each leaf cell of the current SPTDS into two. From the resulting two volumes, two new cells are generated and appended to the SPTDS.

In k-d tree construction, the axis lying orthogonal to the cut plane is chosen periodically, meaning all cells with an equal distance (depth) to the root of the SPTDS are cut with planes orthogonal to the same axis. In the case of ORB, the axis where a cell's volume is largest is picked to compute an orthogonal plane. Therefore cells with the same depth can have cut planes of variable orientations. This makes ORB special as it combines a balanced distribution with near cubic shapes.

2.3 Rise of the GPU

Most if not all state of the art supercomputers are hybrid, meaning each node has a CPU as well as a GPU [TOP500, 2022]. GPU accelerated code can increase the speed of many computations, where in most cases CUDA is used to fully leverage the hardware.

The computational effort required in PKDGrav and most modern astrophysical N-Body simulations can be divided into three categories:

1. **Global Tree Building / Load Balancing:** In supercomputers, particles are distributed equally with regards to a combination of memory and workload among computing nodes to leverage node level parallelization. ORB is used to generate a well suited SPTDS for FMM on a node level.
2. **Local Tree Building:** The same approach is repeated on a local level, where instead of distributing the workload among nodes, its distributed among threads. The tree from step 1. is expanded by the locally computed tree.
3. **Force Calculation and Integration:** The STPDS constructed in step 1. and 2. combined with FMM is used to compute force integration.

In PKDG each category is about one third of the total calculation time [Stadel, 2001]. Making Global (1.) and Local Tree Building (2.) subject to possible performance improvements, as GPU acceleration is as of now only exploited in Force Calculation (3.).

The main objective is to improve the runtime of astrophysical simulations, or more specifically PKDGrav. Simultaneously penalties with regards to N , the total number of particles, should be kept as low as possible. Simulation accuracy benefits from large N as well as an increased number of time steps. As GPU accelerations are low level, hardware specific details are relevant. Thus the algorithms should be adapted and optimized with regards to the target hardware Piz Daint.

2.4 Related Work

A number of GPU accelerated k-d tree construction algorithms exist, however the methods cannot simply be translated to ORB due to different axis choosing methods. For example in the master thesis of Voglsam [Voglsam, 2013] a k-d tree is computed on the GPU and used to improve ray tracing for 3D graphics. However it uses a method to find an optimal cut across cells of the same depth, which does not translate when they have to be found in a variable axis.

ORB can be subdivided into components which have been implemented in existing libraries. Thrust by *NVIDIA* [NVIDIA, 2022g] has an implementation of a reduction and a partition, both are crucial components of ORB. However the library is very high level and does not expose memory management, which is crucial to a highly efficient implementation of ORB. Furthermore CUB [NVIDIA, 2022b], also developed by *NVIDIA*, exposes a more low level API of the components found in ORB but the implementation is not general enough to deal with some critical performance issues. As the SPTDS grows, so do the number of leaf cells and subsequently the number of cells for which a cut plane has to be found simultaneously. When initializing a kernel for each cell, the number of kernels initialized grow exponentially, resulting in kernel invocation overheads dominating the actual computational costs. The problem is almost equivalent to a *Device Segmented Reduce* [NVIDIA, 2022c] where CUB provides an API, however the reduction operator is set to be the same across all segments, which is not true in this case as for each cell the particle coordinates are compared to a different value.

3 Orthogonal Recursive Bisection (ORB)

Each cell in the SPTDS is enriched with the following properties:

- V_{cell} 3D volume where all corresponding particles are contained within.
- d_{cell} Number of leaf cells to be found in the subtree of the cell.
- N_{cell} Number of particles encompassed in the volume of the cell.

Naturally each cell has two pointers pointing to its two children, which are used to traverse the tree, parent cell pointers important but could be added as well.

3.1 Target Data and Notation

A modified random Gaussian distribution is used to generate the target dataset, where it is set to fit measured constraints of the actual universe. Thus the number of particles can be set dynamically.

- The space space of binary numbers with a precision of p is defined as \mathbb{B}_p where $a \in \mathbb{B}_p \Leftrightarrow a \in \{0, 1\}^p$
- The corner coordinates of the root volume are defined as \vec{lower}, \vec{upper} where $\vec{lower}, \vec{upper} \in \mathbb{B}_p^3$.
- The coordinates of a particle i are defined as \vec{x}_i for which holds $\{\vec{x} | \vec{lower} \leq \vec{x} \leq \vec{upper}, \vec{x} \in \mathbb{B}_p^3\}$. When referring to a single coordinate of a particle object, it is done by writing $\vec{x}_{i,j}$. Finally the array of all particle positions in a single axis is denoted as $\vec{x}_{:,j}$.
- N defines the total number of particles in the entire data set. In simulations the number is usually larger than one billion.
- d defines the number of leaf cells the finished SPTDS should have. In the case of load balancing it can be assumed to be no larger than 4096 as most supercomputers do not exceed such a range.

It lies in the nature of the universe, that large clusters of particles are found in some places, whereas other areas may be vastly empty of any objects. This is a very crucial characterization and differentiates the target data set from many other applications of FMM, where the distribution is more uniform.

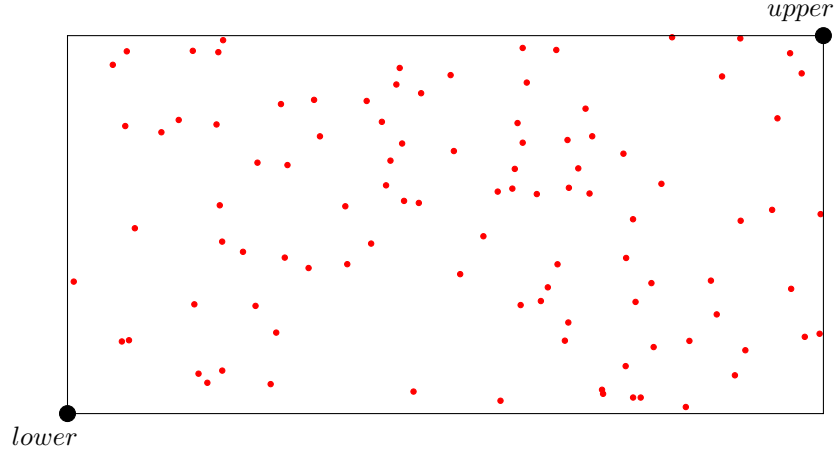


Figure 1: Uniform random distribution of 3D coordinates in cube domain projected onto a 2d plane

3.2 Memory and Workload Balancing

The aim of this paper is to improve the runtime of ORB. In terms load balancing, there is a distinction between workload and memory balancing.

To reduce the force integration error across all particles, some particles require more timesteps than others. This is caused by the vastly variant forces which are exerted on the particles. Whilst some follow a straight path with an almost constant velocity requiring little computational effort for great precision, others are influenced by strong gravitational poles resulting in highly curved movement paths and correlating with higher numbers of timesteps required.

The workload for a particle p_i is defined as the weighting function $w(p_i)$. The workload correlates with the number of simulation steps that need to be computed for a single particle. In an optimally parallelized system balanced in terms of computational effort, the workload should be very similar among computing units. Ideally each core processes a set of particles, such that all cores finish nearly instantaneously. In this case, there is little to no waiting time when synchronizing the results with all involved parallelized cores. With strongly varying weights among particles, perfectly balancing the computational effort yields an unbalanced distribution in terms of memory load.

The degree to which work or memory load balancing is targeted, can be chosen by adjusting the variance in values of the weighting function. Assuming workload balancing should be totally ignored, the weighting function is just set to a constant value for all particles such as $w(p_i) = 1$.

3.3 ORB

To introduce the ORB algorithm, a recursive implementation of ORB is considered as it is easier to understand than the iterative version. All cell properties are defined recursively.

As the recursion is initialized with the root cell alone, it follows that V_{root} is equivalent to the entire volume of the input data. $d_{cell} = d$ and $N_{root} = N$ are true per definition of N and d . When a cell is cut into two child cells, the children are referred to as *leftCell* and *rightCell*. For any cell the following equations hold:

$$depth_{cell} = \lceil \log_2 d_{cell} \rceil \quad (1)$$

$$V_{cell} = V_{leftCell} \cup V_{rightCell} \quad (2)$$

$$V_{leftCell} \cap V_{rightCell} = \emptyset \quad (3)$$

Defining d_{cell} for each cell is especially important, since its crucial to ensure the final tree is a nearly complete binary tree. The way $d_{leftCell}$ and $d_{rightCell}$ is defined, guarantees that on the last level of the tree the cells are filled from the left to the right.

This way no gaps are present and the tree can be stored as a heap in an array.

$$d_{leftCell} = \min\{d_{cell} - 2^{depth_{cell}-2}, 2^{depth_{cell}-1}\} \quad (4)$$

Equation 4 can be explained as follows: When d is a power of two, a perfectly balanced tree can be generated and naturally $d_{leftCell}$ and $d_{rightCell}$ are equal to $\frac{d}{2}$ or $2^{cell_{depth}-1}$. If d is incremented by one so is $depth$ due to the ceiling operator. To have a tree with one more leaf, a cell which has been a leaf node before needs to be divided into two leaf cells and becomes an intermediate cell. The only appropriate choice for this leaf cell is the one on the very left, otherwise the result would be an unbalanced tree. The number of leaf cells in the subtree of *leftCell* is subsequently incremented by one as well, resulting in a number of cells equal to $2^{cell_{depth}-2} + 1$ which is equal to $d - 2^{cell_{depth}-2}$. The equation holds until all leaf cells in left subtree are divided into two child leaf cells and converted into intermediate cells. As this corresponds to a doubling of the number of leaf cells, it can be calculated with $2^{depth_{cell}-1}$. Consequently we limit the number of particles by using a min function and setting $2^{depth_{cell}-1}$ as the second parameter. After $2^{depth_{cell}-1}$ is reached, the subtree of *rightCell* is enlarged until a d corresponding to the next power of two is reached.

$$d_{rightCell} = d_{cell} - d_{leftCell} \quad (5)$$

Equation 5 for $d_{rightCell}$ can also be derived but the more simple variant is just taking the difference between the absolute number of desired leaf cells and the ones that are assigned to the left subtree.

Finally the number of particles encompassed in $V_{leftCell}$ and $V_{rightCell}$ are defined as follows:

$$N_{leftCell} = \min \left\{ x \in \{0, \dots, N_{cell}\} : \sum_{i=0}^x w(p_i) \geq \frac{d_{leftCell}}{d_{cell}} \times \sum_{i=0}^{N_{cell}} w(p_i) \right\} - 1 \quad (6)$$

$$N_{rightCell} = N_{cell} - N_{leftCell} \quad (7)$$

To simplify the visual and numeric explanation of the ORB algorithm, its assumed $\forall i \in \{0, \dots, N\} : w(p_i) = 1$.

The cut plane, which divides $V_{leftCell}$ and $V_{rightCell}$ is axis aligned and thus orthogonal to a given axis a . Thus searching for a single value c is sufficient to construct the plane, where c is the position of the plane along the axis a . The plane is considered ideal if the number of particles where $x_{i,a} \leq c$ are equivalent to $N_{leftCell}$.

After a plane is found, the V_{cell} can be divided or cut in two volumes, where $V_{leftCell}$ and $V_{rightCell}$ are constrained by the original V_{cell} and the cut plane. As defined in equation 3 the volumes do not intersect, thus $V_{leftCell}$ along axis a ends at c , where $V_{rightCell}$ start at c .

3.3.1 Algorithm

Algorithm 1 as mentioned below, is the main routine of ORB. All the volumes are described using a lower and an upper boundary point, sufficiently describing a box, thus the only input parameters required is x storing the particles positions, *lower*, *upper* and finally d . If d is equal to 1, this means the reduction must not be continued, as the target of a single subsequent leaf cell is already reached, since the cell itself is treated as a leaf cell. The stopping condition is formulated on lines 2-4. Line 6 describes a call to a method called *maxIndex()*, which essentially compares all provided values and returns the index of the maximum values. The result i provides us with the axis where the cell volume is largest. d_{depth} and $d_{leftCell}$ are computed as described in equation 1, 4 and the results are stored l and d' (line 7-8). All prerequisites are met to compute the actual cut with the *cut* method which will be explained in more detail in section 3.5 (line 9). The return value stored in *cut* is equivalent to the position of the cut plane

along the cut axis. Using the *cut* value, the array of particles is partitioned as described in 3.6 returning the *mid* value which is the pivot index of the partitioning (line 10). Finally the original volume can be subdivided into two volumes as described on line 11-13 and the *ORB* method is called recursively (lines 15-16) passing it two slices of the *x* array.

Algorithm 1 The ORB main routine

```

1: procedure ORB( $x, \vec{lower}, \vec{upper}, d$ )
2:   if  $d = 1$  then
3:     return ▷ Stopping condition
4:   end if
5:    $\vec{size} = \vec{upper} - \vec{lower}$ 
6:    $i = \text{maxIndex}(\vec{size}_0, \vec{size}_1, \vec{size}_2)$  ▷ Get index of max value

7:    $l = \lceil \log_2 d \rceil$  ▷ Depth
8:    $d' = \min\{d - 2^{l-2}, 2^{l-1}\}$ 
9:    $cut = \text{cut}(x, \vec{lower}_i, \vec{upper}_i, i, \frac{d'}{d})$  ▷ Find cut plane
10:   $mid = \text{partition}(x, split, axis)$  ▷ Partition particles

11:   $\vec{upper}' = \vec{upper}$ 
12:   $\vec{upperChild}_i = cut$ 
13:   $\vec{lower}' = \vec{lower}$ 
14:   $\vec{lowerChild}_i = cut$ 

15:  ORB( $x_{0:mid}, \vec{lower}, \vec{upper}', d'$ )
16:  ORB( $x_{mid:len(x)}, \vec{lower}', \vec{upper}, d - d'$ )
17: end procedure

```

3.4 By Example

A sample dataset is proposed to help with visual and numerical explanations of the algorithm. All particles in the sample dataset are assumed to have a weight of 1.

The algorithm is explained visually using the following example data set:

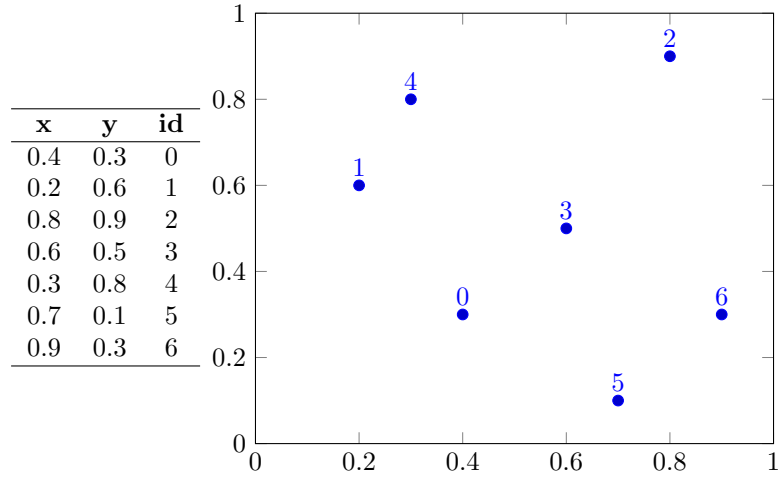


Figure 2: Example 2D particle distribution with $N = 7$

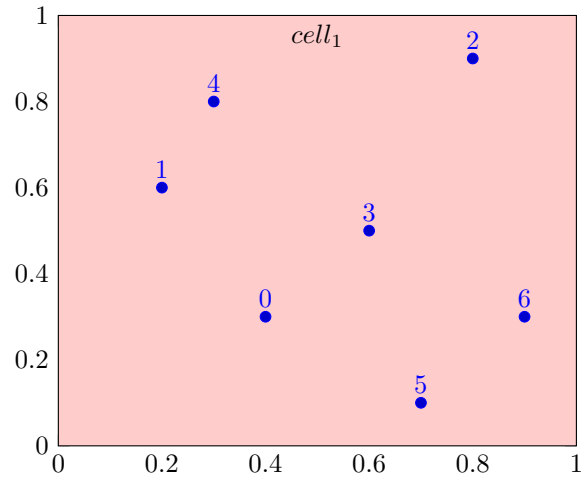


Figure 3: Volumes of SPTDS leaf cells generated from ORB at a depth of 0

$cell_1$
|

Figure 4: SPTDS generated from ORB at a recursion depth of 0

Depicted in figure 3 & ?? is the SPTDS after initialization: there is only one

$cell_1$ which is root and V_{cell_1} encompasses the entire domain which in this case is the rectangle ranging from 0 to 1 in both the x and y axis. d_{cell_1} is chosen to be equivalent to 3 and N_{cell_1} is 7. Note that d could be any other number, but clearly it is limited by N .

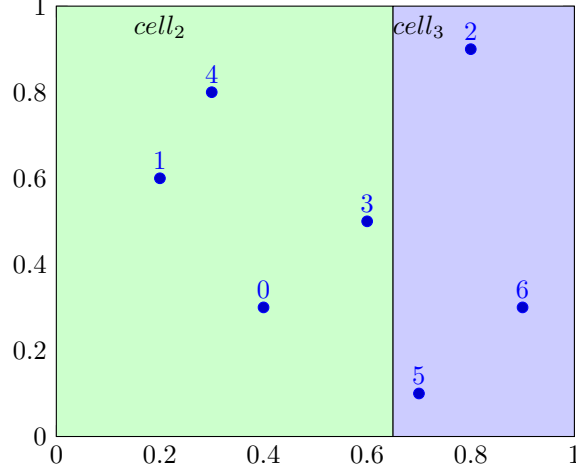


Figure 5: Volumes of SPTDS leaf cells generated from ORB at a depth of 1

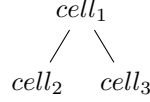


Figure 6: SPTDS generated from ORB at a recursion depth of 1

A cut plane, which translates to a cut line in the simplified 2D example is constructed and $cell_2$ and $cell_3$ are generated accordingly.

5 & 6 Depict a SPTDS of depth 2, or two levels, where a cut plane was found dividing V_{cell_1} into V_{cell_2} and V_{cell_3} , $d_{cell_2} = 2$ and $d_{cell_3} = 1$ subsequently are calculated, meaning $cell_3$ is considered a leaf cell. Thus N_{cell_2} can be computed using a simplified version of equation 6 since all weights are assumed to be 1: $N_{cell_2} = \lfloor \frac{d_{cell_2}}{d_{cell_3}} * N_{cell_1} \rfloor = \lfloor \frac{2}{3} * 7 \rfloor = 4$ subsequently $N_{cell_3} = 3$.

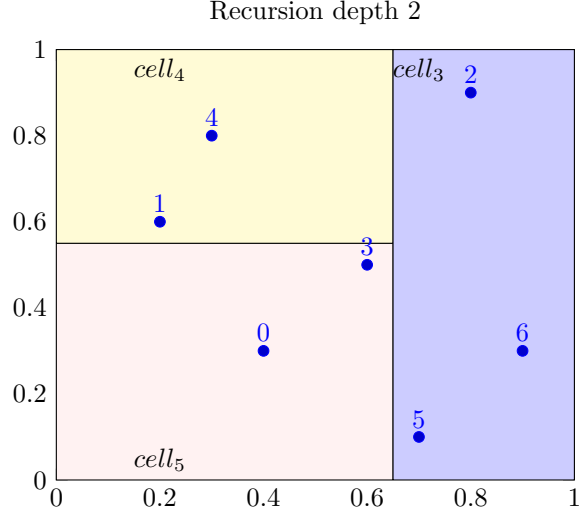


Figure 7: Volumes of SPTDS leaf cells generated from ORB at a depth of 2

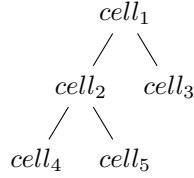


Figure 8: SPTDS generated from ORB at a recursion depth of 2

Again the procedure is repeated for $cell_2$ but not for $cell_3$ as there recursion is terminated by the stopping condition. Finally $cell_4$ and $cell_5$ are computed, where the stopping condition is met as well. The resulting SPTDS has 3 leaf cells and partitions the space into three subspaces.

3.5 Root finding

The *cut* algorithm takes an array of particle positions x , an *axis*, *left* and *right* boundaries and a *percentage*. Its goal is to return a position along the cut axis such that the particles less or equal to the a *cut* value are equivalent to the percentage multiplied by the length of the x array.

The problem is related to a selection algorithm, as selection algorithms are capable of finding the k th smallest element in an array of numbers. To adapt it for ORB one can set $k = f \times N_{cell}$, where f is the desired percentage of particles which should be assigned to the *leftCell* of *cell*. A quick select

[Cormen et al., 2022] could be used, however quick select has a worst case runtime of $O(n^2)$. Furthermore median algorithms could be explored as well, but they are mostly approximation algorithms with bad worst case runtimes. The problem can be reformulated as a root finding problem. To do so a function $f(c)$ is defined which evaluates the number of elements smaller than the cut value minus half¹ the number of total particles. When the function evaluates to zero, an axis aligned cut plane is found, where exactly half the particles are located on the left side. There exist many different solvers for the root-finding problem, but the most stable and easiest to implement is the bisection method. Some solver combine approximate solver with the stable bisection method to generate fast but stable root finding methods. Such algorithms could be explored in later work.

Since the algorithm operates on binary numbers with a limited precision, the bisection method is guaranteed to terminate after p steps.

3.5.1 Bisection Method

Initially an estimation for a cut is made, in this case the exact middle of the domain boundaries. Some median finding algorithms use improved guessing to speed up the process, again a method which could be leveraged in later work.

Because the maximal number of iterations is known, a loop can be used as seen below in algorithm 2 on line 4. A cut is then computed (line 5), which in the first iteration is the center position between left and right. It is then checked whether the stopping condition has already been reached (line 7-9), meaning the cut lies within α points of the ideal result. If it was not reached, the boundaries can be improved as follows: In case there are too many elements left of the cut, it is known that the cut position was chosen too far to the right and it follows that the ideal cut must be left of the current guess. Therefore, the boundaries can be adjusted, in this case by setting *right* equal to *cut* as seen on line 13. The analogous concept can be applied in the other case (line 11).

¹Or any other desired percentage

Algorithm 2 Bisection Method

```
1: procedure CUT( $x, left, right, axis, percentage$ )
2:    $nLeft = 0$ 
3:    $cut = 0$ 

4:   for  $k \in 0, \dots, p$  do
5:      $cut = (right + left)/2$ 
6:      $nLeft \leftarrow \text{sum}(x_{:,axis} < cut)$  ▷ Counting particles left of cut

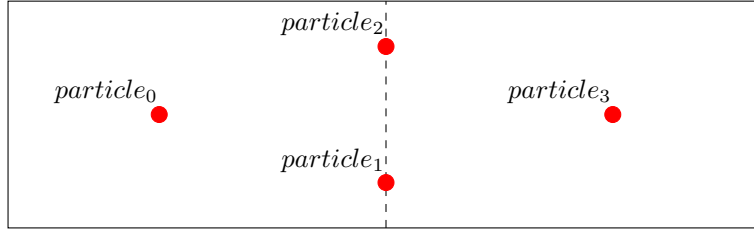
7:     if  $\text{abs}(nLeft - \text{len}(x)) * percentage < \alpha$  then
8:       Break ▷ Stopping condition
9:     end if

10:    if  $nLeft \leq \text{len}(x) * percentage$  then
11:       $left = cut$ 
12:    else
13:       $right = cut$ 
14:    end if

15:  end for
16:  return  $cut$ 
17: end procedure
```

3.5.2 Edge Cases

Considering the following example particle distribution where $particle_1$ and $particle_2$ have identical x coordinates.



In this case there exists no ideal cut in the x axis². As either there is 1 particle to the left or 3, but no cut can result in 2 particles to the left. If more particles are added along the same line, the method performs even worse. However since the algorithm uses a deterministic for loop, the algorithm will terminate anyways and puts either 3 particles to the left or the right.

²x axis is horizontal, the y is vertical

3.5.3 Runtime Analysis

On line 5, the number of particles to the left of the cut plane are summed. As the array is unordered, all coordinates of one axis need to be read once, resulting in a runtime of $O(N)$. All other operations inside the for loop can be computed in a negligible constant time. The loop itself is repeated p number of times, where p corresponds to the precision of x coordinates.

The algorithm must terminate after p iterations for integer numbers. With each iteration the range of possible solutions is divided by two. The same goes for integer numbers, each time a bit is removed, the size of the range of numbers which can be represented is reduced by two. Thus after p iterations the precision limit is reached and the cut cannot be improved. The proof does not hold analogously for float numbers as the precision is uniformly distributed there, however its also limited.

3.6 Partition Algorithm

Its desired to continuously update the array storing the positions of the particles and in the words used before, have a direct correlation between x_i and i . This enables grouping particles within a cell in a fixed range of two indexes of the particles array. The advantages of this are two fold: Access all particles within a cell in constant time, manipulate particles within a cell using a slice of the array. Another option is to store the assigned leaf cell for each particle, this however requires a much larger array of size N .

As seen in the pseudo code in algorithm 3, the algorithm looks for a pair of particles, where for both the coordinate along the relevant axis are on the wrong side of the cut plane. In this case the particles can be swapped (line 8) resulting in the correct position of both particles. Correctness and a more detailed explanation of the algorithm can be found in the book [Cormen et al., 2022].

Algorithm 3 Partition Method

```

1: procedure PARTITION( $x, cut, axis$ )
2:    $i = 0$ 
3:   for  $k \in 0, ..N - 1$  do
4:     if  $\vec{x}_{k,axis} \leq cut$  then
5:       while  $\vec{x}_{i,axis} \leq cut$  and  $i < N$  do
6:          $i = i + 1$ 
7:       end while
8:        $x_i, x_k = x_k, x_i$ 
9:     end if
10:  end for
11:   $x_i, x_{N_j-1} = x_{N_j-1}, x_i$ 
12: end procedure

```

A runtime of $O(N)$ can be derived, as the algorithm iterates over all particles once. Since each element in the array needs to be touched at least once to

partition the entire array there exists no better method.

Applying the algorithm to the running example looks as follows:

x	y	id
0.4	0.3	0
0.2	0.6	1
0.8	0.9	2
0.6	0.5	3
0.3	0.8	4
0.7	0.1	5
0.9	0.3	6

The particles are then partitioned with a cut in the x axis set to 0.65 as seen in figure 5.

x	y	id
0.4	0.3	0
0.2	0.6	1
0.3	0.8	4
0.6	0.5	3
0.8	0.9	2
0.7	0.1	5
0.9	0.3	6

Finally *cell*₂ is cut divided into *cell*₄ and *cell*₅ with the cut position 0.55 along the y axis as seen in figure 5 ending up with:

x	y	id
0.4	0.3	0
0.6	0.5	3
0.3	0.8	4
0.2	0.6	1
0.8	0.9	2
0.7	0.1	5
0.9	0.3	6

Note how all particles from *cell*₄ are contained in the range of 0-1. The particles of *cell*₅ in 2-3 and finally the particles contained in the volume of *cell*₃ can be found in rows 4-7.

4 Theoretical Analysis

The main goal of this paper, is to improve the runtime of the ORB algorithm by leveraging the graphics processing unit over the central processing unit. Before implementing, it makes sense to explore if and to which degree the performance benefits could in theory manifest themselves. For this purpose a simplified model is proposed, this model is then used to compare the runtime of the CPU version against two different GPU variants. As GPU programming is very low level, the knowledge gained while developing the model also builds a solid theoretical foundation.

4.1 General Memory Model

For a consistent terminology of a computer, a general computing model, which can be applied to most modern high performance systems is established. The model reflects a single node and its hardware components, a supercomputer may have thousands of these nodes linked together where different bandwidths are seen between individual nodes. However the internodal communication and parallelization is not important in this work, because the communication and optimization strategies are handled by parts of the PKDGrav which we adopt for this research. Both the CPU and the GPU have their own memory which are connected by a data link, where the connections are bound by the memory bandwidths. The capacity of the CPU memory bandwidth is named B_{CPU} and the GPU memory bandwidth B_{GPU} . There is a separate data link between the CPU and the GPU memory, which is commonly referred to as PCI express or NVLink (for modern *NVIDIA* GPU's, here the simplified term I_{GC} is used. In systems with multiple CPU's there is a link between the individual CPU's which is denoted as I_{CC} . Finally the data link between GPU's is referred to as I_{GG} .

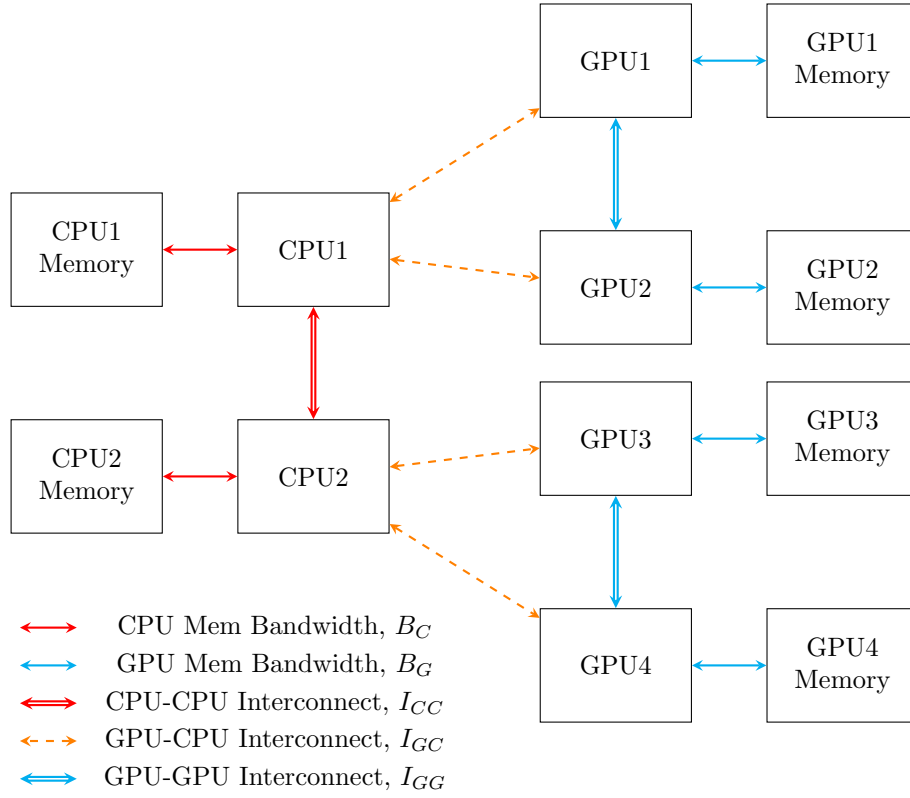


Figure 9: Computing model

4.2 Supercomputers

For Piz Daint, Summit and Eiger all relevant hardware metrics were collected and compiled to table 1. Piz Daint and Eiger were chosen, because the systems can be used to test the code. Furthermore Summit is included as it is, at the time of writing this paper, one of the most capable supercomputers in the world.

Constant	Piz Daint [CSCS, 2022]	Summit [ORNL, 2022]	Alps (Eiger)
# Nodes	5704	4608	1024
# CPUs	1	2	2
CPU Model	Intel E5-2690 v3 [Intel, 2022]	IBM POWER9	AMD EPYC 7742[AMD, 2022]
CPU Mem.	64 GB	256 GB	??
B_C	68 GB/s	170 GB/s	204.8 GB/s x 2
I_{CC}	-	64 GB/s	??
Base GHZ_C	2.9 GHZ	4 GHZ	2.25 GHZ
Max GHZ_C	3.8 GHZ	4 GHZ	3.4 GHZ
# Cores	12	22	64
Architecture	Haswell	POWER9	AMD Infinity Architecture
SIMD	AVX2	VSX	AVX2
# GPUs	1	6	0
GPU Model	NVIDIA P100 [NVIDIA, 2022f]	NVIDIA V100s [?]	-
GPU Mem. Cap.	16 GB	16 GB x 6	-
B_G	732 GB/s	900 GB/s x 6	-
I_{GC}	32 GB/s	50 GB/s x 6	-
I_{GG}	-	50 GB/s	-
GPU Tflops	9.3	16.5	-
# CUDA Cores	3584	5120	-

Table 1: Datapoints of Supercomputers

4.3 Roofline Performance Model

In a first step it is determined whether the computations are bound by memory bandwidth or the actual performance of the computing unit. This can be done using a roofline model. It compares arithmetic intensity in flops per byte against the actual performance of the computing chip in flops and can provide insights whether an application is memory or performance bound.

The most costly computation is line 6 of the bisection method (figure 2) with a runtime of $O(32 \times N)$. Oftentimes when an algorithm iterates over a large dataset performing only very little calculations on its individual elements, the limiting factor is the memory. To support this claim, a roofline model for all three systems is established.

4.3.1 Estimating Flops

For modern hardware, it is fairly uncommon to release flops (floating point operations per second) values. Steadily evolving SIMD instruction sets result in varying performance for different implementation details, which in turn are compiled into different assembly instructions. Depending on the algorithm, implementation details and compilation flags the c++ compiler tries to compile ideal assembly instruction sets. SIMD instructions can only be used when there is a contiguous memory access, thus in some cases a poor memory layout choice may lead to a much lower flops.

For most modern CPU chip architectures AVX is the fastest SIMD instruction set available. The common AVX2 enables the processing of 8 floating point operations per instruction with a cpi^3 of 0.5. The cpi varies depending on the chip architecture, however all relevant CPU's from table 1 support a CPI of 0.5 along with AVX2. A lower cpi results in a higher efficiency, as several instructions can be completed in a single cycle.

Equation 8 defines a function to estimate the number of gigaflops for a given hardware. hz is the number of cycles which can be completed within a single second

³cycles per instruction

by a processing unit, meanwhile nf is the number of floats which can be processed simultaneously in a single instruction using SIMD. Finally np represents the number of processors, where a perfect parallelization is assumed, meaning 100% of the code can be parallelized.

Variable	Meaning	Unit
s	second	
cl	processor cycle	
f	floating point operation	
i	instruction	
(g)hz	(10^{12}) processor cycles per second	(10^{12}) cl / s
nf	floating point operations per instruction	f / i
cpi	cycles per instruction	cl / i
(g)flops	(10^{12}) floating point operations per second	(10^{12}) f / s
np	number of processors	

Table 2: Variables to estimate gflops

$$gflops = ghz * nf * cpi^{-1} * np \quad (8)$$

Since peak flops benchmarks are available for *NVIDIA* P100 and V100s no estimations need to be made on the GPU side.

4.3.2 Estimating Arithmetic Intensity

Arithmetic intensity is measured in FLOPS per byte or the number of floating point operations which are computed per byte loaded from memory. The count left part from the cut algorithm (line 6) in figure 2 can be translated to the following isolated c++ code:

Listing 1: Counting the particles left of a cut plane

```
1 for (auto p= startPtr; p<endPtr; ++p) nLeft += *p < cut;
```

Where p is a C-style array which stores the particle's position and $nLeft$ stores the number of particles which are smaller than cut . The operations per loaded float can be listed as follows:

1. Compare particle to cut
2. Add result to $nLeft$
3. Increment pointer p
4. Compare pointer p with $endPtr$

Which results in a total of 4 operations. Since a single float is stored using 4 bytes, this computes to $\frac{4flops}{4bytes} = 1.0$ operations per byte or an arithmetic intensity of one.

Estimating the arithmetic intensity for the GPU is a lot more complicated as it can vary a lot depending on the specific implementation details. But for now the same arithmetic intensity is assumed for both CPU and GPU.

Note that SIMD instructions do not influence the arithmetic intensity, as the number of floating point operations per byte remains the same. It just means that a set of operations is executed concurrently. What is however influenced by AVX, is the maximum number of GFLOPS which can be processed by the hardware. If AVX was ignored, the algorithm would clearly not be bound by memory, but by the much lower performance of the processing unit. For this reason it is important to consider SIMD instructions.

4.3.3 The Plots

The maximally achievable gflops for Piz Daint are computed as follows: The datapoints from figure 1 are plugged into the formula 8 for $np = 2, 4$ and 8.

$$3.8 \times 10^{12} cl/s \times 8f/i \times \frac{1}{0.5} i/cl \times 2 = 121.6 gflops \quad (9)$$

$$3.8 \times 10^{12} cl/s \times 8f/i \times \frac{1}{0.5} i/cl \times 4 = 243.2 gflops \quad (10)$$

$$3.8 \times 10^{12} cl/s \times 8f/i \times \frac{1}{0.5} i/cl \times 8 = 486.4 gflops \quad (11)$$

The results are depicted as horizontal lines in figure 10. The memory limit is plotted as a line with a slope equivalent to its respective bandwidth. Finally a dotted line represents the arithmetic intensity of the Count Left procedure. To interpret the roofline mode, one has to follow the dotted line starting from the bottom. Whether it intersects with the line representing the maximal performance or the memory bandwidth first, gives an indication whether the program is performance or memory bound.

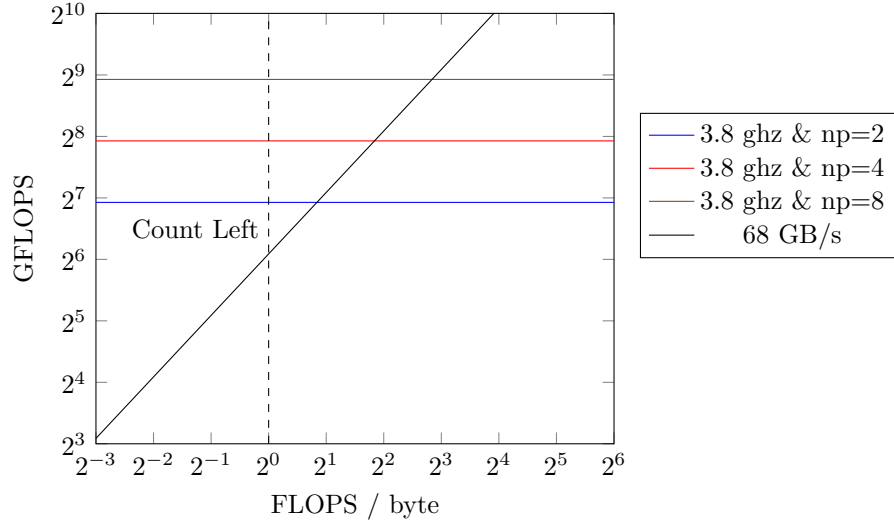


Figure 10: Roofline Model for the isolated Count Left code on Piz Daint

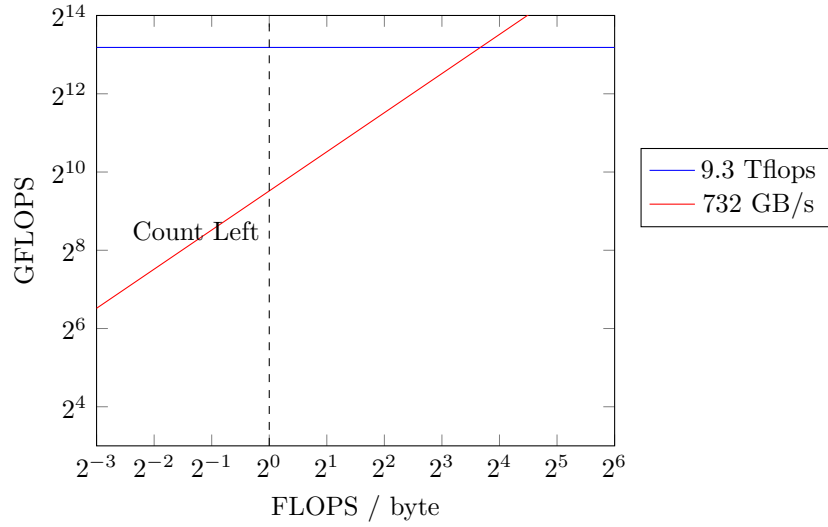


Figure 11: Roofline Model for the isolated Count Left code on Piz Daint using the GPU

As seen in figure 10 and 11 both the GPU and CPU versions running on Piz Daint are memory bound. When tracing the vertical line representing the arithmetic intensity of the Count Left procedure, it crosses the memory limit before the processing

unit limit. However the memory limit on the GPU is much higher, therefore it should be favored over the CPU.

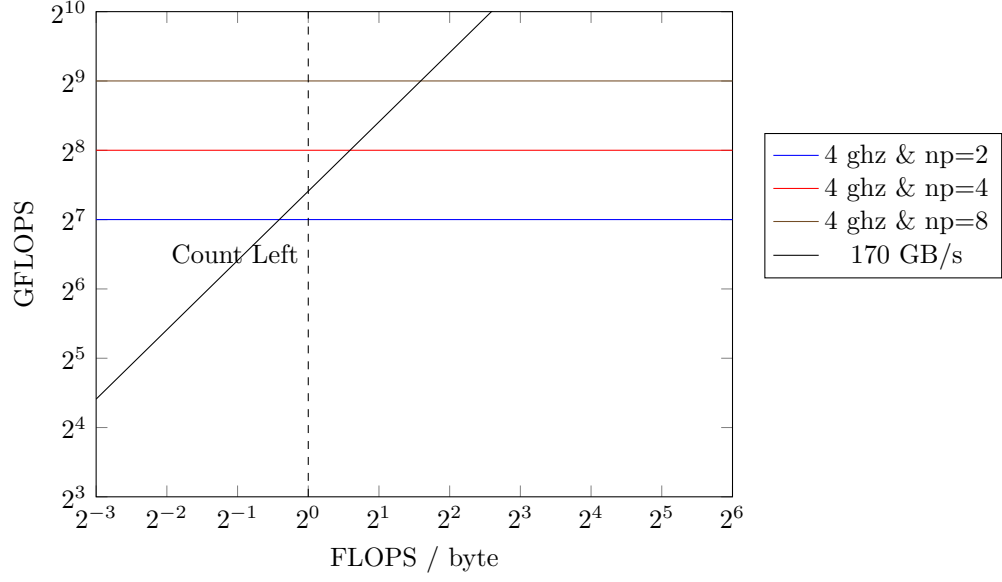


Figure 12: Roofline Model for the isolated Count Left code on Piz Daint

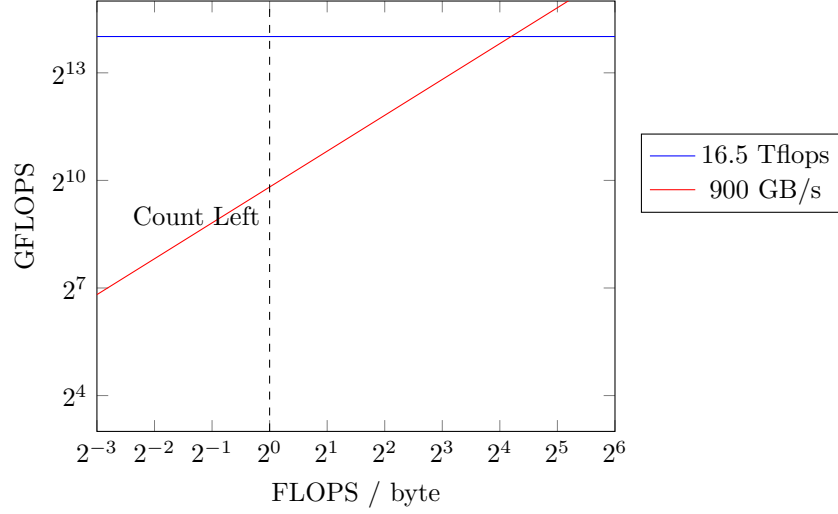


Figure 13: Roofline Model for the isolated Count Left code on Summit using the GPU

As visible in figure 12 and 13 the same applies to Summit. Its notable how the memory bandwidth only becomes the limiting factor when assuming a parallelization with 4 processors.

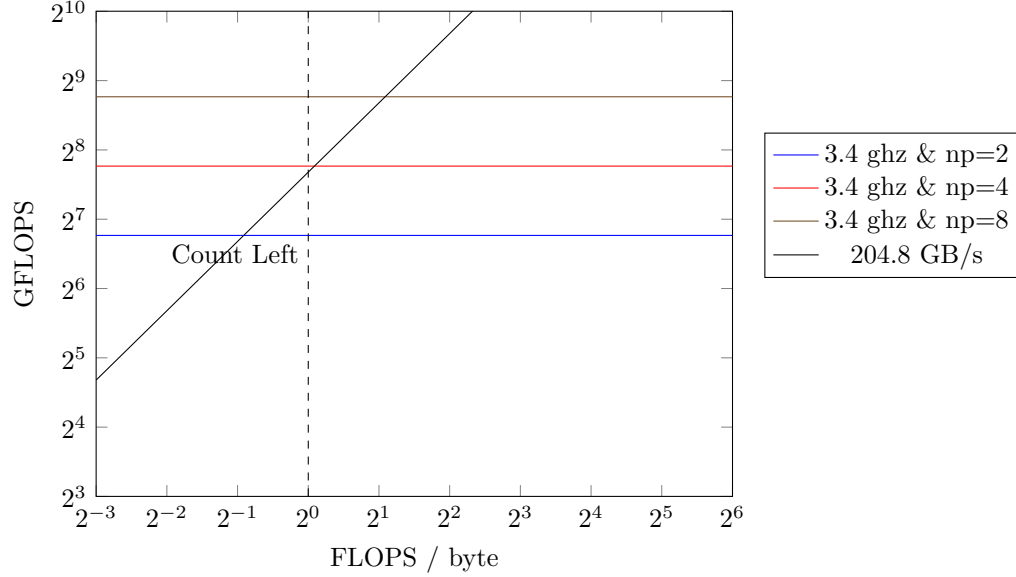


Figure 14: Roofline Model for the isolated Count Left code on Alps

Finally alps, which has the highest memory bandwidth compared to its flops, also becomes limited by the memory bandwidth after using more than four processors.

4.3.4 Empirical Verification

A minimal C++ code is considered to verify the theoretical model. The -S flag along with the g++ compiler is used to generate assembly code from the c++ source code. For testing purposes all entries in the array are set to random values between 0 and 1, and *cut* is set to 0.5.

In a first test, only O3 is enabled. The generated assembly code is depicted in listing 2.

Listing 2: Reduction Assembler Code without AVX

```

1 .L18:
2     movups    (%rax), %xmm0
3     addq     $16, %rax
4     cmpltps   %xmm2, %xmm0
5     psubd    %xmm0, %xmm1
6     cmpq     %rdx, %rax
7     jne      .L18

```

vpsubd is a packaged instruction, meaning it already uses some form of SIMD instructions. The command is used in the MME and later the SSE2 instruction set. Since the *xmm* registers are used, it follows that it is an AVX instruction.

If the compile flag *-march = native* is added, the instructions as depicted in listing 3 are the result of the compilation.

Listing 3: Reduction Assembler Code with AVX2

```

1 .L19:
2     vmovups  (%rax), %ymm3
3     addq    $32, %rax
4     vcmpltps %ymm2, %ymm3, %ymm0
5     vpsubd   %ymm0, %ymm1, %ymm1
6     cmpq    %rdx, %rax
7     jne     .L19

```

vmovups, *vcmpltps* and *vpsubd* are AVX2 commands since they are using the *ymm* registers.

The following results were obtained on 2^{27} particles with a thread count of 16 on Piz Daint. The reduction takes 2.3 milliseconds, resulting in a throughput of $2^{27}/10^9 * 10^3/2.3 = 58gflops$. This is close to the memory bandwidth (68 GB/s) and far below the performance maximum of the CPU, which even for ($np = 1$) is $3.8 * 8 * 2 * 1 * 16 = 972.8gflops$. It concludes that the algorithm is from a theoretical and empirical point of view memory bound on the GPU. The evidence for the GPU is not as clear as the arithmetic density cannot be estimated as well, however when assuming a similar density as the CPU variant, the algorithm is also on the GPU bound by memory.

4.4 Runtime Estimates

A runtime estimation for the CPU version and two different GPU implementations is constructed. In the first GPU implementation *GPU Counting*, an array of particles is sent to the GPU, where counting is then performed. In a more advanced variant *GPU Counting and Partitioning* the particles are sent from the CPU to GPU only once and the partitioning is performed on the GPU as well.

A precision p of 32 is assumed, which is a sensible assumption for astrophysical simulations. The total storage required for all particle positions is $32 \times 3 \times Nbits = 4 \times 3 \times NBytes = 12 \times NBytes$. Furthermore $d = 4096$ and $N = 10^9$ are assumed. The sizes are realistic, considering Piz Daint has a GPU memory limited to 16 GB and 5704 nodes.

Variable	Meaning	Unit
d	number of leafs cells in the final SPTDS	
p	precision of each coordinate (number of bits)	
s	data size of all particle coordinates	
B_C	memory bandwidth between the CPU and its memory	GB/s
B_G	memory bandwidth between the GPU and its memory	GB/s
B_G	memory bandwidth between the CPU and the GPU	GB/s

Table 3: Variables used for runtime estimates

4.4.1 CPU Version

Each time the leaf cells of the SPTDS are split into two child cells and appended to the tree, the tree ends up with twice the number of cells. As the goal is to construct a tree with d leaf cells, the cut algorithm needs to be performed $\lceil \log_2(d) \rceil$ times over all leaf cells.

Consequently the runtime estimate sums over all $\lceil \log_2(d) \rceil$ iterations, where in each iteration a cut has to be found for i cells. The cut for a single cell is found in p iterations, where the number of particles to be iterated over for a single cell is divided by two in each iteration.

$$\sum_{i=1}^{\lceil \log_2 d \rceil} i \times p \times \frac{s}{B_C} \quad (12)$$

i cancels out from the equation and the equation becomes:

$$\sum_{i=1}^{\lceil \log_2 d \rceil} p \times \frac{s}{B_C} \quad (13)$$

Which can be simplified to:

$$\lceil \log_2(d) \rceil \times \left(p \times \frac{s}{B_C} \right) \quad (14)$$

4.4.2 GPU Counting

The Equation 15 for the GPU is similar, the only difference being that GPU memory bandwidth B_{GPU} is used instead of the CPU bandwidth. Furthermore the additional overhead of sending the data from the CPU to the GPU has to be reflected. This adds the term size divided by CPU to GPU memory bandwidth denoted as I_{GC} .

$$\lceil \log(d) \rceil \times \left(p \times \frac{s}{B_G} + \frac{s}{I_{GC}} \right) \quad (15)$$

4.4.3 GPU Counting and Partitioning

Another variant is to perform the partitioning on the GPU, meaning there is no need to send data from the CPU to the GPU each time a cut needs to be found. This allows for a reduction of the costly overheads introduced by the rather slow bandwidth I_{GC} . The overhead can be moved out of the bracket as the transfer needs to happen only once before.

$$\lceil \log(d) \rceil \times \left(p \times \frac{s}{B_G} \right) + \frac{s}{I_{GC}} = t \quad (16)$$

4.4.4 Plugin Values

Lets plugin the values from figure 1 into the corresponding formulas 14, 15 and 16 for Piz Daint. The following values and consequently speed up can be computed for CPU only version:

$$\lceil \log(4096) \rceil \times \left(32 \times \frac{12GB}{68GB/s} \right) = 67.77406s \quad (17)$$

Whereas a GPU accelerated Count Left results in the following runtime:

$$\lceil \log(4096) \rceil \times \left(32 \times \frac{12GB}{732GB/s} + \frac{12GB}{32GB/s} \right) = 10.79648s \quad (18)$$

This yields in a speed-up of the GPU over the CPU only version:

$$\frac{67.77406}{10.79648} = 6.2774 \times \quad (19)$$

When considering the GPU Counting and Partitioning the results are:

$$\lceil \log(4096) \rceil \times \left(32 \times \frac{12GB}{732GB/s} \right) + \frac{12GB}{32GB/s} = 6.67085s \quad (20)$$

This yields in a speed-up over the CPU obly version of:

$$\frac{67.77406}{6.67085} = 10.15973 \times \quad (21)$$

4.5 Conclusion

The computations for Eiger and Summit can be performed analogously to the ones from Piz Daint. Eiger is non hybrid, thus only estimates for a CPU version are computed. All the results are plotted in 15. As expected the GPU Counting outperforms the CPU version on both hybrid⁴ supercomputers. Furthermore GPU Counting and Partitioning yields more performance improvements, but the speed up is less drastic. As the model may differ from reality depending on the implementation and compilation details, these results cannot be taken for granted but offer a good baseline of what is possible.

⁴Each node has a CPU and a GPU

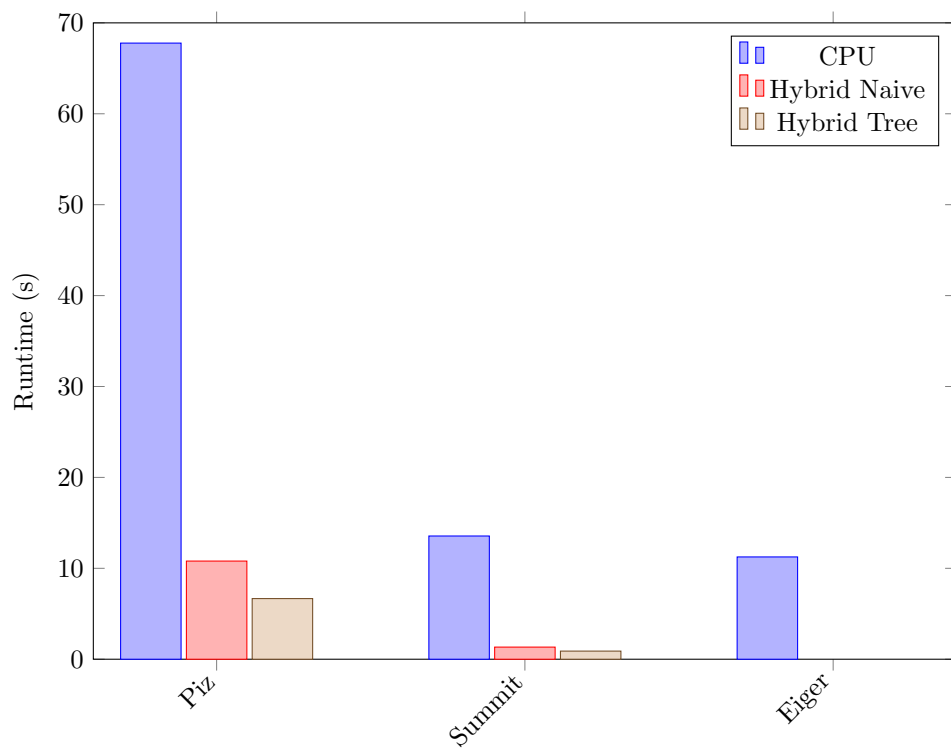


Figure 15: Estimated runtimes of different variants

5 CPU Implementation

In this section the stand-alone implementation of the ORB algorithm using CPU acceleration only is described. Building on the CPU version, GPU accelerated alternatives to selected parts of the code are proposed.

The code is developed somewhat separately from the PKDGrav codebase in order to keep the project smaller and more easy to understand. However the machine dependent layer (MDL) from PKDGrav is used, which provides a system to distribute the workload among processors. This reduces the complexity of the implementation because communication and CPU parallelism are abstracted away from the code. Furthermore it makes integration into PKDGrav simple. Not every part of the code is described in detail, however any details crucial to performance are outlined. Furthermore, special attention is paid to the CUDA kernels and memory management strategies for two reasons: minor implementation details can have a big impact upon the performance and increased replicability.

For this project C++ 17 along with CUDA 11.3, OpenMPI 4.1.0 and Blitz++ 1.0.2 [Veldhuizen, 2022] is used.

5.1 MDL and PKDGrav

MDL, the machine dependent layer, provides an abstraction layer around the parallel primitives introduced by PKDGrav. The Master layer is responsible for coordinating the flow of the program. It does so by calling the PST (Processor Set Tree), which distributes the tasks and gets processors to work on them. As the name suggests, the PST is organized as a binary tree where intermediate nodes contain a pointer to a consecutive processor as well as a next lower node. Each leaf cell of the graph then correspond to an individual core. [Stadel, 2001]

Parallel processes are dispatched by descending the PST until they reach all processors. There computations are performed and the results are combined by passing them back up the PST until they reach the master. A parallel process can be programmed by implementing the service interface as explained in section 5.5.

Each processor has its own local data, which can be accessed by calling `pst->lcl`.

5.2 Particles

The particles array has ownership over all particle objects. It has a (N,3) shape whereas each row represents a single particle. To enable coalesced memory access when iterating over a single axis i.e. $x_{:,axis}$, a column major storage order is chosen, meaning that each column is stored in a range of consecutive array addresses. Even though row access patterns i.e. $x_{i,:}$ are used as well, they happen less frequently. Most importantly there are cases where entire columns have to be copied, which is by a magnitudes faster with a column major storage.

C-style arrays are used interchangeably with Blitz++ arrays to store all particles. Blitz++ is an open source wrapper class around C-style arrays, which helps with pointer management and can speed up the debugging process by keeping track of the array boundaries. Furthermore, it features array slicing and does its own memory management. In certain cases the actual C-style array needs to be accessed, which can be accessed by calling `array.data()`. For example when iterating over the Blitz++ array, non optimal assembly code was observed, not reflecting optimal SIMD instructions available to the hardware.

The particle data is loaded into the local data storage of each processor, meaning each processor owns a different unique set of particles.

5.3 Cell

The cell class is a structure keeping track of the fundamental cell information. In essence it is the analogue to the concept of the *cell* which has been introduced in section 3.

Quickly building and accessing elements of the SPTDS generated by ORB requires a suitable data structure. Instead of a tree with pointers, a heap can be used since the necessary conditions for a nearly complete binary tree are met. Elements within a heap can be accessed and added in constant runtime. Depicted in figure 16 is the finished SPTDS of the sample dataset stored as a heap.

The SPTDS and thus the cells are constructed on the master alone. Cells are only distributed when dispatching services on the PST.

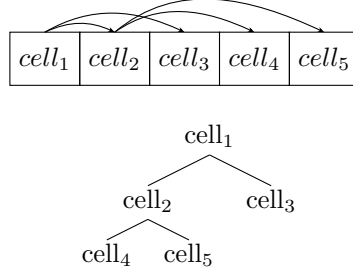


Figure 16: Tree stored as an array heap

An array can be used as an efficient storage medium for the heap, since in this case d , the number of leaf cells in the final SPTDS, can be set at compile time and the maximum number of cells in the SPTDS can be derived from that. Meaning the array can be allocated statically. Because MDL communicates data as arrays between threads, another advantage of the array storage is the absence of any costly data structure conversion.

5.3.1 Heap Conditions

As with all heap data-structures the following conditions are given:

- The root cell has an index 1
- The left child of a cell at index i can be accessed at the index $i \times 2$.
- The right child of a cell at index i can be accessed at the index $i \times 2 + 1$.

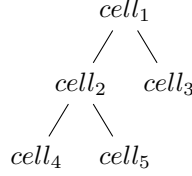


Figure 17: SPTDS constructed from example particle distribution tree with 3 leaf cells ($d = 3$)

The following constraints can be derived given the number of leaf cells d .

- The depth of the tree can be computed with: $\lceil \log_2(d) \rceil$
- The number of leaf cells on the second last level is given by $2^{depth} - d$
- There are exactly $2 \times d - 2^{depth}$ items (which must all be leaf cells) on the last level.
- The total number of cells are $2 \times d - 1$.

In the example tree depicted in figure 17 the depth corresponds to: $\lceil \log_2(d) \rceil = 2$, the number of leaves on the second last level is $2^2 - 3 = 1$ and the number of leaves on the last level are $2 \times 3 - 2^2 = 2$. Finally the total number of cells is $2 \times 3 - 1 = 5$.

5.3.2 Class

Among the fields in the cell are the following variables:

- *id* is the unique identification of the cell instance. It corresponds to its index in the heap array plus one. The plus one comes from the different indexing for heap constraints, starting with 1, and the classic array indexing starting with 0. It can be used to compute indices of both child cells and parent cells using the above shown formulas.
- *nLeafCells* is equivalent to d_{cell} and depicts the number of leaf cells to be found in all of its successors. It is used when building the tree, to track whether a cell needs to be split or not.
- *lower* is a 3D point coordinate representing the lower boundary corner of the 3D volume V_{cell} which is encompassed by this cells.
- *upper* represents the upper boundary corner of the volume.
- *cutAxis* is the axis where the cell is to be cut.

To support and simplify common operations, a further abstraction layer is given with the following functions:

- *setCutAxis()* evaluates the proper *cutAxis* given the *upper* and *lower* variables.
- *cut()* cuts the cell and returns two new cells.
- *getLeftChildId()* returns the ID of the left child cell, meaning its volume is in the cut axis smaller than the cut plane.
- *getRightChildId()* is the equivalent for the right child.

5.4 Mapping Cells to Particles

As the SPTDS heap is stored on the master but the particles are distributed among the local storages of the processors, it is crucial to know which particles are encompassed in the volume of which cell. Thus a data structure to keep track of the relation between cells and its particles is introduced. Thanks to the partitioning algorithm all particles belonging to a single cell can be found in a consecutive section of the particles array. A 2D array is used to keep track of the relation, which again is stored in the local data of the processor. For each cell the corresponding section is stored as a tuple of begin and end indices. Naturally the data structure needs to be updated across all processors whenever a partition is performed.

5.5 Services

A service implementation consists of a *Service()* function as well as a *Combine()* function. Parameters include an *PST* object which stores information about the process location as well as the local data. Furthermore two void pointers providing storage for the input and output data called *in* and *out* along with their respective sizes are passed. The reason void pointers are used is the variability of the input data structures. C++ does not offer dynamic typing. Data stored inside the void pointers can be converted into their respective classes using casting.

In the *Combine()* function two output void pointers *vout* and *vout2* and the sizes of the underlying arrays are given as parameters. A combination strategy of the elements can be chosen and implemented, where the results is to be stored in *vout*. For example in a service which returns the sum of all particles, the combine function would simply store the sum of *vout* and *vout2* in *vout*.

5.5.1 Init and Finalize Service

All data stored in the local memory of a PST, must be allocated when initializing the program. After the program has finished, the memory is freed to avoid any leaks. Allocation and freeing comes with some overhead, thus it makes sense to reuse memory whenever possible.

The initialize service is also responsible to load the particle data into the local memory of each PST.

5.5.2 Count Left Service

As mentioned, counting the elements smaller than the cut position is computationally intensive. The *Service()* function which performs the actual computations, *in* points to an array of cells which is given by a void pointer and as the length of the array is known, its possible to simply iterate over all of its elements and cast them to the correct cell data structure as shown in listing 4 on line 1-2. In case the *foundCut* flag of the cell was set to true, it is known that the ideal cut was already found and the cell can be skipped (line 4-6). Then the begin and end indices of the particles array corresponding to the range of objects contained within the cells volume are read and respectively a slice of the local particles array (lines 7-11) is made. Start and end pointers are computed and the cut position to be tested is fetched (lines 13-17). Finally the number of particles smaller than the cut are counted and written back to the output array (lines 18-23). As the results for each cell needs to be stored individually, its written to the output with an offset equivalent to the *cellPtrOffset*.

This offset however does not correspond to the *cell.id*, its rather a sequential indexing over all leaf cells of the same depth in the SPTDS.

Listing 4: Part of the Count Left Service() method

```

1 for (int cellPtrOffset=0; cellPtrOffset<nCells; ++cellPtrOffset){
2     auto cell = static_cast<Cell>*(in + cellPtrOffset);
3
4     if (cell.foundCut) {
5         continue;
6     }
7     int beginInd = pst->lcl->cellToRangeMap(cell.id, 0);
8     int endInd = pst->lcl->cellToRangeMap(cell.id, 1);
9
10    blitz::Array<float,1> particles =
11    pst->lcl->particlesAxis( blitz::Range(beginInd, endInd));
12
13    float * startPtr = particles.data();
14    float * endPtr = startPtr + (endInd - beginInd);
15
16    int nLeft = 0;
17    float cut = cell.getCut();
18    for(auto p= startPtr; p<endPtr; ++p)
19    {
20        nLeft += *p < cut;
21    }
22
23    out[cellPtrOffset] = nLeft;
24 }
```

Combining the data is straightforward, as in *Combine()* *vout* and *vout2* contain the counts for a given number of particles. The master requires a global count across all the particles contained within all local storages, therefore each element of *vout* is summed together with *vout2*.

Listing 5: Part of the Count Left Combine() method

```

1 for(auto i=0; i<nCounts; ++i)
2     out[i] += out2[i];
```

5.5.3 Count

Counting the total number of particles contained within a cell is trivial with the *cellToRangeMap*. The results can be obtained by subtracting the begin index from the end index of each cell. The *Combine()* is equal to the one shown for the Count Left Service.

Listing 6: Part of Count Service() method

```

1 for (int cellPtrOffset=0; cellPtrOffset<nCells; ++cellPtrOffset) {
2     auto cell = static_cast<Cell>*(in + cellPtrOffset);
3     out[cellPtrOffset] =
4     lcl->cellToRangeMap(cell.id,1) -
5     lcl->cellToRangeMap(cell.id,0);
6 }
```

5.5.4 Partition

The Partitioning Service is a direct implementation of the Hoare partition and can be seen in listing 7. Again the service receives a buffer with a given number of cells, which are then iterated over using the outer for loop. Each cell from the buffer is then casted to the proper struct (line 2). Using the id of the cell, the start and end indices of the corresponding particles stored in local storage can be accessed (lines 4-5). Lines 7-37 search for pairs located on the wrong side of the cut and then swap them using a simple *swap()* helper function. *i* has to be initialized using *beginInd* - 1 as the *i* ++ operation happens before a read on the array (line 15). For the same reason *j* is initialized to *endInd* and not *endInd* - 1. Then the while loop terminates, *i* points to the first index in the array where the value is greater than the cut, therefore *cellToRangeMap* can be updated using *i*, as for the left child cell the start remains equal but the end is now *i* and for the right child cell its vice versa. Note that the implementation is not ideal for two reasons:

1. The blitz wrapper class should be avoided and instead the memory directly accessed.
2. In this particular case the column major storage order is not ideal, as rows are accessed over and over again.

Listing 7: Partition Service

```

1 for (int cellPtrOffset=0; cellPtrOffset<nCells; ++cellPtrOffset){
2     auto cell = static_cast<Cell>*(in + cellPtrOffset);
3
4     int beginInd = pst->lcl->cellToRangeMap( cell.id , 0);
5     int endInd = pst->lcl->cellToRangeMap( cell.id , 1);
6
7     int i = beginInd-1, j = endInd;
8     float cut = cell.getCut();
9
10    while(true)
11    {
12        do
13        {
14            i++;
15        } while( lcl->particles(i, cell.cutAxis)<cut and i<=endInd);
16
17        do
18        {
19            j--;
20        } while( lcl->particles(j, cell.cutAxis)>cut and j>=beginInd);
21
22        if(i >= j) {
23            break;
24        }
25
26        swap(lcl->particles , i , j);
27    }
28
29    swap(lcl->particles , i , endInd -1);
30
31    lcl->cellToRangeMap( cell.getLeftChildId() , 0) =
32    lcl->cellToRangeMap( cell.id , 0);
33    lcl->cellToRangeMap( cell.getLeftChildId() , 1) = i;
34
35    lcl->cellToRangeMap( cell.getRightChildId() , 0) = i;
36    lcl->cellToRangeMap( cell.getRightChildId() , 1) =
37    lcl->cellToRangeMap( cell.id , 1);
38
39 }

```

5.5.5 Make Axis

As mentioned in section 2.1 the relevant axis to search for the cut position differs for each cell. Therefore the algorithm cannot simply iterate over a single array. In order to simplify the process the Make Axis service is introduced. For each cell it copies a single axis of particle coordinates encompassed in the respective volumes. Consequently the axis is equivalent to the cut axis, meaning the axis where a cut plane is to be found.

5.6 Parallel Schedule

In the context of parallelization, the number of processors is defined as np and p_0 is the processor executing the master layer. Initially each processor has a random unique

subset of all the particles stored in its local memory.
 In the running example this might look as follows:

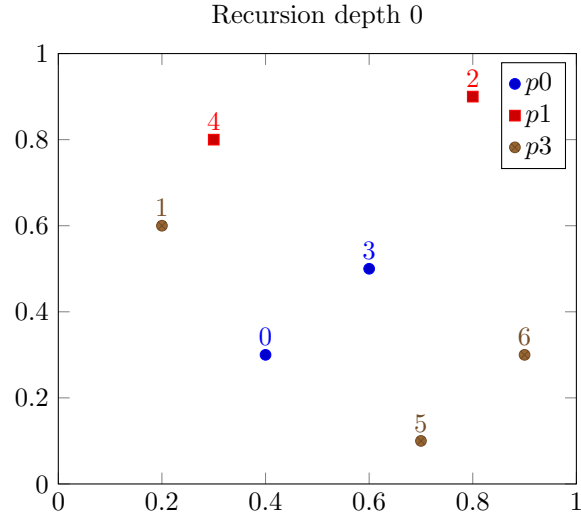


Figure 18: Example particles distributed randomly across 3 nodes

The Master schedules all other processors but also performs tasks by itself. The exact schedule is illustrated in figure 19 where a service which is dispatched from the master and executed on all processors is depicted as a horizontal rectangle. Computations which are only performed by a single processor are represented by a vertical rectangle.

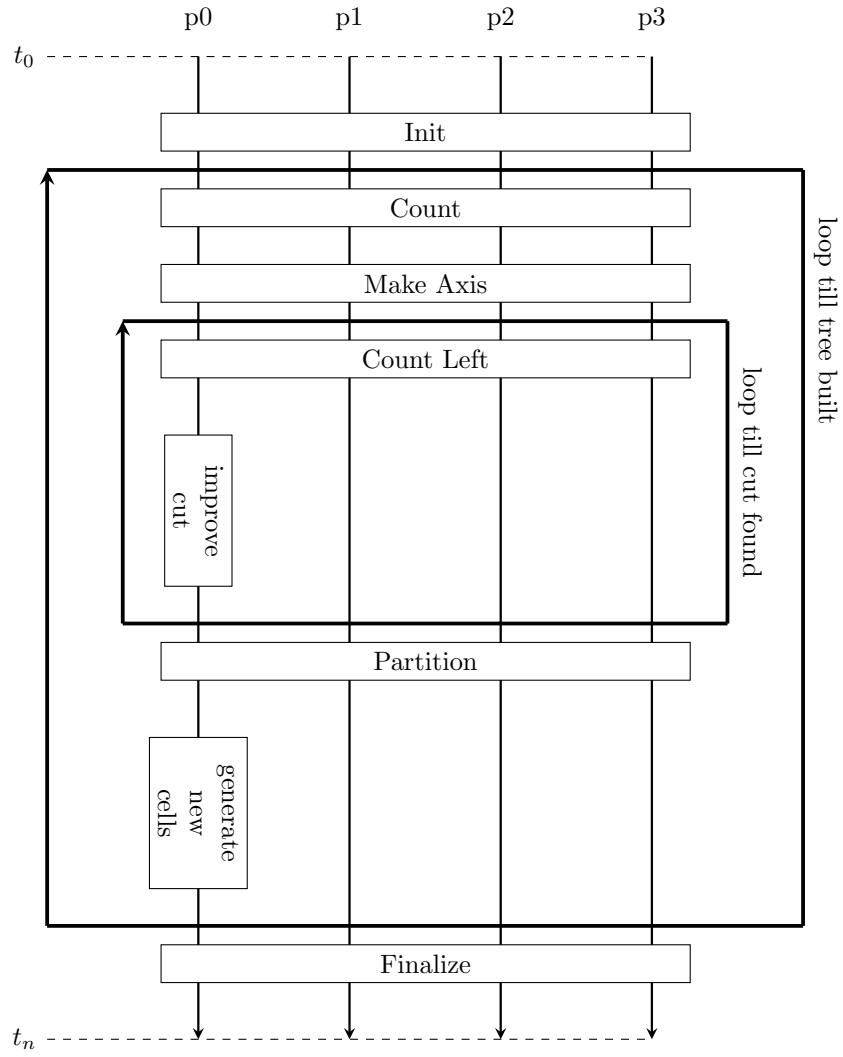


Figure 19: Master later schedule of a CPU only ORB implementation

6 GPU Implementation

CUDA is a parallel computing platform and programming model developed by *NVIDIA*. [NVIDIA, 2022e] It allows software developers and programmers to use a CUDA-enabled graphics processing unit (GPU) for general purpose processing. The CUDA platform is designed to work with programming languages such as C, C++, and Fortran and gives direct access to the GPU's virtual instruction set and parallel computational elements.

6.1 Relevant CUDA Concepts

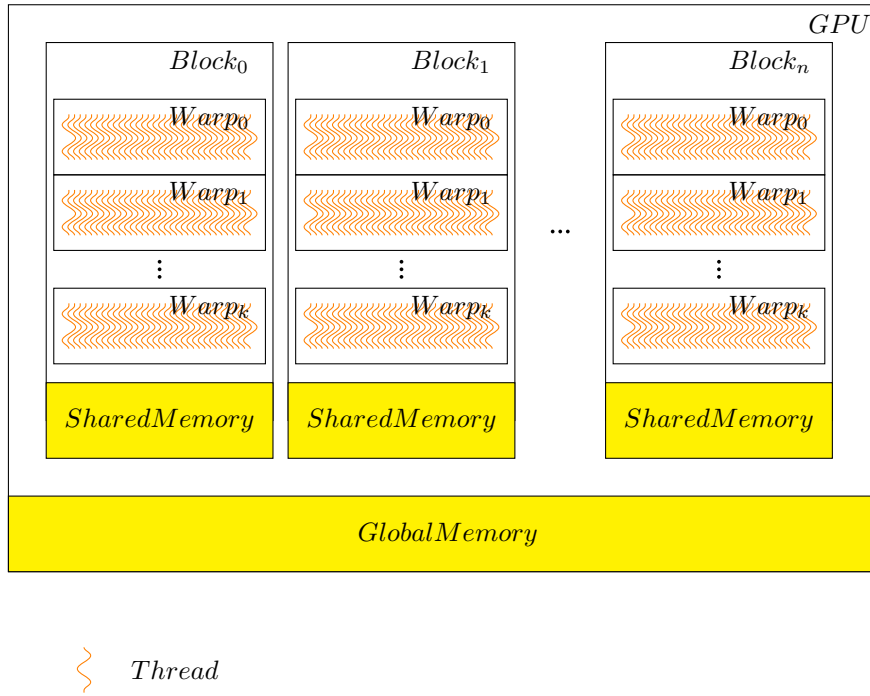


Figure 20: Fundamental CUDA concepts

CUDA gives the ability to launch kernels, which are written in C with some additional CUDA specific syntax. These kernels can be run from the CPU, commonly referred to as the host, and are then executed on the GPU, also known as the device. The device and host have a separate memory and as seen in section 4.1 the transfer I_{CG} rates between the GPU and CPU are generally slower, under performing I_C and I_G speeds. Therefore one of the key challenges when rewriting CPU code to GPU code is to limit data transfers between the GPU and CPU as much as possible. Furthermore the problems have to be divided into smaller subproblems, such that each of them fits a single block. CUDA cannot give any guarantees considering the order of execution of these blocks, greatly limiting the possibilities of communication between blocks.

Algorithms have to be fundamentally redesigned when porting them from CPU to GPU code. Generally problems which can be vectorized, meaning the same computation is applied to every element in an array are well suited. Usually vectorization is accomplished by mapping each thread associated with a kernel to a value in the array where it performs the respective calculations. As communication is costly and cross block communication is strongly discouraged, computation where threads are less dependent on results from other threads are better suited. Finally problems with less control statements work better, as control statements can cause significant overhead.

The location of a thread within a block can be accessed using *threadIdx.x*. The block ID among all other blocks associated with the same kernel is stored in *blockIdx.x*. Finally the number of blocks are stored in the variable *gridDim.x*. [NVIDIA, 2022e]

Building a tree is therefore a rather challenging problem as there are many computations which influence another and building a tree involves a lot of branching, at least when it is done in a conventional manner.

Listing 8: Calling a CUDA Kernel

```

1 add<<<
2     nBlocks ,
3     nThreads ,
4     nThreads * bytesPerThread ,
5     stream
6 >>>(
7     param1 ,
8     param2
9 );

```

A kernel is executed exactly once for every kernel, where a block consists of many kernels. The number of kernels per block and the total number of blocks can be defined by the user as seen in the kernel call syntax in listing 8 on line 2 and 3.

6.1.1 Warps

Each consecutive grouping of 32 threads form a warp. All threads within a warp are executed in parallel, given that there is no warp divergence. A warp divergence can occur, if there are control statements, where two or more threads from the same warp execute a different code. This leads to a decrease in performance and should be avoided whenever possible. [Yuan Lin, 2021]

Additionally warp level primitives can be accessed on a warp level, where they can operate on local instead of shared memory. The primitives enable safe communication between local memory of threads within the same warp, which is not possible otherwise. [NVIDIA, 2022e] [Yuan Lin, 2021]

6.1.2 Memory

Each block has access to a shared memory register, where the capacity of this register can be defined at runtime by the host as seen in figure 8 on line 4. The maximum shared memory size depends on the architecture ranges between 48 and 96kB [NVIDIA, 2022e]. There are also upper limits for the number of threads per block which is usually between 512 and 1024 [NVIDIA, 2022e]. With a shared memory size of 48kb and 1024 threads per block there is $\frac{48000B}{1024} = 46B$ of memory available per

thread, which is roughly equivalent to 11 32 bit precision numbers per thread. However as the thread count per block is usually not maxed out, there is more shared memory space available per thread.

Data in the shared memory exists only as long as the kernel exists and it cannot be accessed across different blocks. Race conditions do apply to shared memory, but threads in a block can be synchronized, therefore it is not safe to have writes from multiple threads on the same shared memory address. There exists a CUDA command `__syncthreads()` which allows us to synchronize all threads within a block, enabling us to control the execution order of individual statements.

Additionally to the shared memory, there exists global memory. Global memory is persistent even after a block or a kernel has finished and is the only way to send data between the device and the host. Since the data is not deleted after a kernel has finished, the data can and should be reused whenever possible. Shared memory latency is roughly 100 times lower than global memory latency. [Harris, 2021b] and it is best practice to copy data from global memory to shared memory before performing actual computations on it. Reducing the total number of memory accesses on global memory. After the computations have finished, the results can be written back from shared to global memory. [Harris, 2020a]

6.1.3 Memory Access Patterns

Global memory is fetched from memory in 32 Byte packages, which translates to $\frac{32}{4} = 8$ floats. A simple aligned sequential memory access pattern, where each thread reads a single float from global memory, results in 4 transactions per warp (32 threads). Unaligned memory access can have minor performance penalties because more memory banks have to be loaded per warp. [Harris, 2020a]

Each 32th entry in shared memory is stored in the same memory bank. If from the same half warp, multiple threads access the same memory bank in shared memory, it results in a bank conflict. If 16 consecutive threads access 16 different banks, then there is no bank conflict and all the data can be fetched in parallel.

6.1.4 Asynchronous Operations

There are two types of engines that can be used to execute kernels in CUDA streams: copy engines and kernel engines. Copy engines are used to copy data between host and device memory, and between different types of memory on the device. Kernel engines are used to execute CUDA kernels. The number of individual engines depends on the actual hardware, lower end GPUs usually have a single kernel and a copy engine, more advanced architectures can have more than one copy or kernel engine.

A CUDA stream is a sequence of commands that are executed in order on a CUDA device [Rennich, 2021]. Streams can be used to improve the runtime of a CUDA program by overlapping the execution of different kernels. For example, a copy kernel can be executed in one stream while a compute kernel is executed in another stream. This overlap can lead to a significant performance improvement.

6.1.5 Synchronization

Whilst threads within a block can be synchronized, this is not the case for the blocks launched from a kernel. The only real synchronization technique is to wait for a kernel,

meaning all blocks associated with this kernel, to finish and execute a consecutive task using another kernel.

Since the GPU is limited in its computing and also memory capabilities, the number of blocks which are run in parallel are limited, therefore CUDA cannot give us any guarantee, whether a set of blocks are run serially or in parallel. Some of the blocks might be run in parallel, meanwhile others are run serially. This very constraint forces the programmer to rethink algorithms and think of ways how individual subproblems can be run somewhat independently of each other.

There exists however a way to communicate between blocks in a safe way. Atomic operations can be performed on global memory without any race conditions.

6.2 Streams

As described in section 6.1.4, streams can be leveraged to improve the runtime of a CUDA accelerated application [NVIDIA, 2022e] [Rennich, 2021]. In the GPU accelerated ORB streams are used in a straightforward manner, each thread performs all its computations using its own unique stream. Since leveraging streams only improves the runtime in low numbers because the actual copy and kernel engines are not very numerous, this variant is very simple and effective. It was observed that using more streams on the individual threads does not improve the runtime, quite contrary it can have a negative impact.

6.3 Memory Management

Several strategies are used to improve the performance of memory allocation:

1. **Pinned Memory** As the GPU cannot access the default paged memory directly, when copying memory from the host to the device, the memory must first be copied to pinned memory. This means when pinned memory is used directly with `cudaMallocHost()`, data transfers can be around twice as fast. Thus pinned memory is used for all data, which needs to be either sent from the host to the device or vice versa [Harris, 2021a].
2. **Reuse Memory** Memory allocation and freeing can be avoided altogether by reusing previously allocated memory whenever possible. Since there exists a fixed number of particles as well as an upper limit for the number of cells, memory can be allocated when calling the Initialize service and finalized with the Finalize service.

6.4 GPU Accelerated Count Left

A reduction problem is fairly similar to the Count Left Procedure. In this case a compare against the cut value is made on all elements, when smaller an int value of 0 results otherwise 1. Then all the results of the compare are summed to a single result. This is a map reduce problem, where the compare maps each input to a desired value, where the values are then reduced across all elements and stored.

6.4.1 Schedule

The entire schedule of the GPU accelerated ORB is depicted in figure 21. In a first step the initialize service is called, where the necessary data is allocated on all devices

and the particles are loaded. Furthermore the initial SPTDS is constructed, which essentially only consists of the root node, encompassing the entire domain and all the particles. Next the main loop of ORB is entered, which iterates until a SPTDS with the desired size is constructed. Within the loop the Count Service is called, computing the number of cells for each cell summed over all the threads in the system. To prepare the data for the GPU transfer a temporary array is made, using the make axis service, which is then sent to the GPU. At this stage the actual bisection method starts on the master, where counting the number of particles left of a cut plane is done with the GPU Count Left service. Depending on the outcome, the cut position is improved respectively and the process is repeated until a nearly perfect cut is found for all the leaf cells of the current SPTDS. For each leaf cell two new child cells can be generated using the computed cut positions. Finally the SPTDS can be enriched with the new cell data and the particles array is partitioned accordingly using the Partition service. The tree building loop is then repeated or if the desired size of the SPTDS is reached, the loop is exited and the Finalize service is called to free the allocated memory.

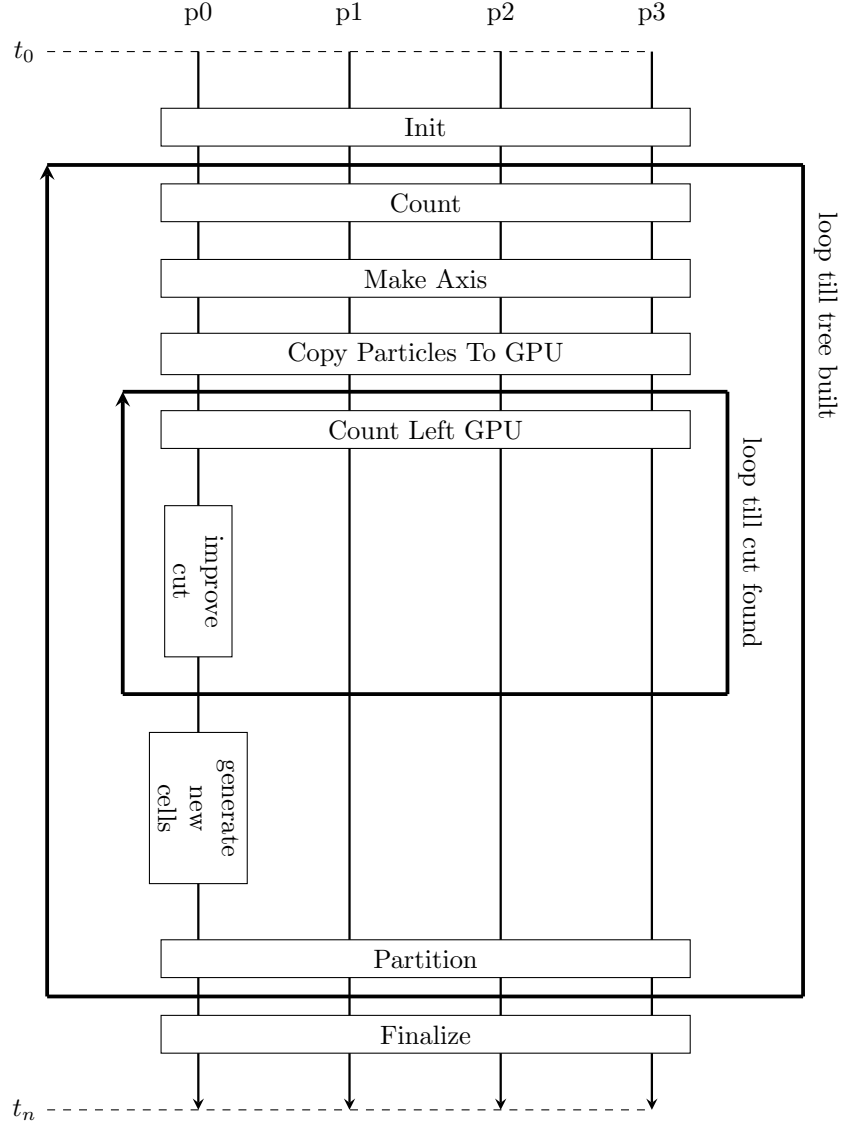


Figure 21: Master later schedule of a GPU accelerated ORB implementation

6.4.2 Service

The GPU version of the Count Left Service invokes the reduction kernel exactly once for every leaf cell of the current SPTDS. Meaning on level 0 of the tree there is only one kernel per processor being executed, when proceeding to further levels, the number of kernels is equivalent to 2^{depth} per processor. The number of blocks per cell can be determined by N_{cell} divided by the number of threads per block and divided by the

number of elements per thread, which is a strategy to increase the thread occupancy.

All the relevant information can be passed along with the cell can be passed using pass by value parameters, there is no need to copy any additional data from the CPU to the GPU in advance. The only relevant data is already in global memory, when the particles were copied from the host to the device. However the results, which are also stored in global memory, need to be copied back from the GPU to the GPU after the kernel has finished.

6.4.3 Kernel Code

The actual kernel code is depicted in listing 9. The input parameters of the reduction kernel (line 17-20) are: a pointer to the the particles stored in global memory, a pointer the an allocated array were the results are stored, the position of the cut plane and finally the total number of elements contained within the cell.

In a first step on lines 21-23 each thread ID *tid* is associated with a pointer offset *i* to the input data *g_idata*. Furthermore on line 25-26 the shared data is initialized and set to zero. Its size is equivalent to the block size, the total number of threads per block. Meaning there is room for a single float per thread in shared memory. On lines 28-30 the thread iterates over all elements which it is associated to and copies the results of a compare to shared memory. This allows us to work exclusively with shared memory instead of global memory, reducing the amount of costly memory accesses.

There exists an optimum in terms of the number of operations performed by a single thread. In order to adapt the number of ops dynamically, the while loop ensures that each thread iterates over a certain number of elements as seen in lines 28-30. The input parameter *n* defines the total number of particles. Therefore if the total number of blocks is reduced by a factor of *r*, the while loop will iterate over *r* elements. Inside the while loop on line 27 the index *i* is incremented by *gridSize* in each iteration. The grid size is equivalent to the total number of threads associated with the kernel. Accessing the global memory in this pattern is coalesced but not necessarily aligned as the pointer of the input array probably does not start at a 32 index aligned address. [Harris, 2020a] [Harris, 2022]

The rest of the kernel code performs the actual reduction, where in each step two elements from the shared memory are summed together. The exact pattern is chosen, because this way shared memory bank conflicts can be avoided, as each thread accesses only elements in shared memory which are stored in the same bank. [Harris, 2021b] Given a block size of 256, after line 39⁵ exactly 256 sums are stored in shared memory. In a next step as seen on line 40 threads 0-127 in the block add the value stored at the same index in shared memory together with the values stored in a position which is offset by 128. This step reduces the total number of elements to be considered to 128, and the step can be repeated as seen on lines 45 to 48 until we reach 64 elements in shared memory which can be reduced by 32 threads. Meaning we only need to perform operations on a warp level which is done in the device kernel called *warpReduce* as it is invoked on line 51.

The *warpReduce* method takes the shared memory as an input and computes a reduction on warp level. Lines 5-6 perform a similar pattern as was seen before, but in this case we only need to synchronize among all threads in the warp using the *_syncwarp()* method. Finally we can perform the reduction on 32 elements by loading them to local memory and using the warp level primitive *_shfl_down_sync(m, v, o)*. The

⁵lines 31-36 are not relevant as the block size is smaller than 512

call gets the value of v of thread with index offset by o compared to the thread which calls the method. The mask m specifies which threads should execute the function. Note that a non full mask allows some form of divergence between threads in a warp without performance penalties. [Yuan Lin, 2021]

The reduction is optimized by using a template which takes as argument the total number of threads per block. Since the template is evaluated at compile time, all the if statements taking `blockSize` as a parameter in figure do not cause any performance loss. In fact, because the reduction loops are unrolled (34-49) the performance of the kernel is increased. [Harris, 2022]

Listing 9: Conditional Reduction in CUDA

```

1 #define FULLMASK 0xffffffff
2 template <unsigned int blockSize>
3 extern __device__ void warpReduce(
4     volatile unsigned int *s_data, unsigned int tid) {
5     if (blockSize >= 64) s_data[tid] += s_data[tid + 32];
6     __syncwarp();
7     if (tid < 32) {
8         unsigned int val = s_data[tid];
9         for (int offset = 16; offset > 0; offset /= 2)
10             val += __shfl_down_sync(FULLMASK, val, offset);
11         s_data[tid] = val;
12     }
13 }
14
15 template <unsigned int blockSize>
16 extern __global__ void reduce(
17     float *g_idata,
18     unsigned int *g_odata,
19     float cut,
20     int n) {
21     unsigned int tid = threadIdx.x;
22     unsigned int i = blockIdx.x*(blockSize) + threadIdx.x;
23     unsigned int gridSize = blockSize*gridDim.x;
24
25     __shared__ int sdata[blockSize];
26     sdata[tid] = 0;
27
28     while (i < n) {
29         sdata[tid] += (g_idata[i] <= cut);
30         i += gridSize;
31     }
32     __syncthreads();
33
34     if (blockSize >= 512) {
35         if (tid < 256) {
36             sdata[tid] += sdata[tid + 256];
37         }
38         __syncthreads();
39     }
40     if (blockSize >= 256) {
41         if (tid < 128) {
42             sdata[tid] += sdata[tid + 128];
43         } __syncthreads();
44     }
45     if (blockSize >= 128) {

```

```

46     if (tid < 64) {
47         sdata[tid] += sdata[tid + 64];
48     } __syncthreads();
49 }
50 if (tid < 32) {
51     warpReduce<blockSize>(sdata, tid);
52 }
53 if (tid == 0) {
54     g_odata[blockIdx.x] = sdata[0];
55 }
56 }

```

Initializing the kernel once for each cell is a simple solution, however as the number of cells is increased, more calls to the kernel are made. Calls are not free of cost, as there is some overhead with each initialization of a kernel. Thus the performance degrades gradually as the tree is built, because on each level the number of cells for which the bisection method is executed doubles. Some performance evaluations have shown that the runtime in the last iteration of ORB, finding a cut for all cells took around 100 times longer than the very first iteration for the root cell.

6.5 Improved GPU Accelerated Count Left

To solve the degrading performance problem some changes are made to the kernel and the overall schedule. The main idea is to prepare the necessary data for all cells in advance, effectively reducing the overall number of necessary kernel calls to $\lceil \log_2 d \rceil$ instead of $2 \times d$. In figure 22 the resulting number of kernel calls are compared for both versions. Each block is provided with information concerning the start and end index of the cell's particles, the cut position as well as an index of the cell.

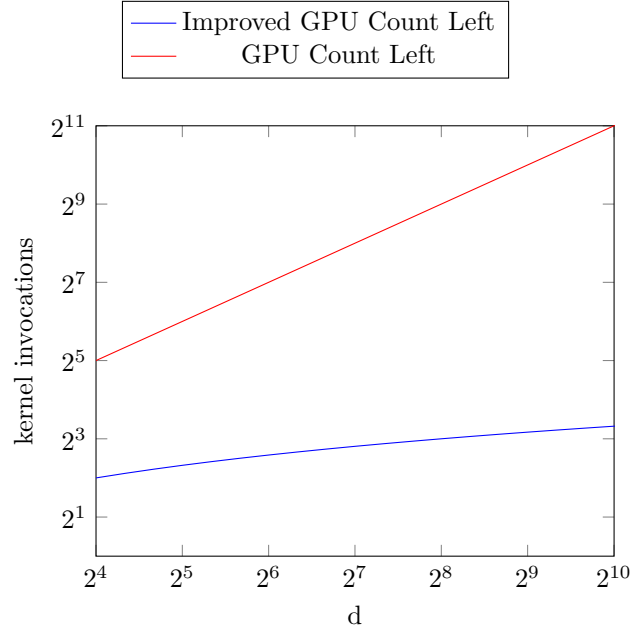


Figure 22: Kernel calls with GPU Count Left with two different versions

The improved Count Left is explained with a simple example: Given $cell_2$ with a volume that encompasses particles in the range 0 - 10240 and $cell_3$ encompasses particles in the range 10240 - 20480. A block size of 256 threads is given and the number of elements per thread is set to 16, thus each block processes $256 * 16 = 4096$ elements. In this case two blocks are allocated to each cell, where the first block processes elements 0 - 4096 and the second one 4096 - 10240. By assigning more work to the second block, each thread in the second block performs more than in the first. However this is a simple solution to avoid blocks, where not even all threads are mapped to an input element, because the number of input elements is smaller than the number of threads. Consequently the third and fourth block respectively process particles 10240 to 14446 and 14446 to 20480. Since the number of required blocks for a cell is known as well as the exact range of particles for which the reduction has to be computed, all this information can be copied to the GPU in advance and a kernel with the total block sum can be launched. Each kernel stores its information in an output array where the assignment between the output and the actual cell is known. The service therefore needs to iterate over all the output which have been computed by the blocks associated with a kernel and output them at a single cell specific location in the output *ou* array of the service.

6.5.1 Schedule

Only a single change is made to the schedule: A service called GPU Copy Cell Service is introduced, which does as the name says, copies the cell information from the CPU to the GPU. Since only the cut positions change while ORB is iterating to find and

improve a cut and the rest of the cell information stays static, this service can be called outside of the loop. The Improved GPU Count Left Service copies the cut data each time it is invoked. Other than that, the schedule remains equal.

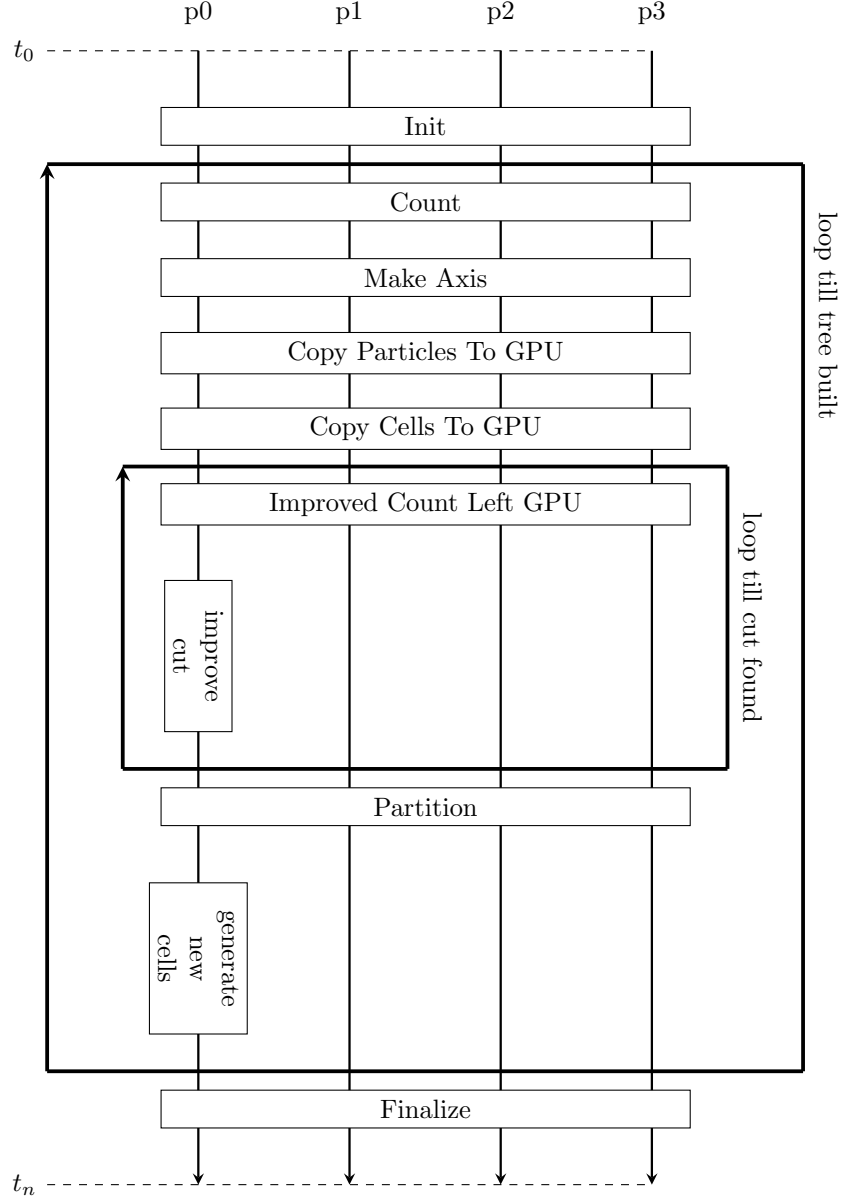


Figure 23: Master later schedule of an improved GPU accelerated ORB implementation

6.5.2 Kernel Code

The kernel code as seen in listing 10 largely remains the same. New data pointers are passed along as a parameter where *g_{begin}* and *g_{end}* mark the begin and end of the important particle array slice, furthermore *g_{cuts}* is an array of cell cuts for each block. The data is prepared in such a way that each block only needs to operate over the particles contained within a single cell. *blockIdx.x* provides the index to read out the cell information from the input arrays (lines 20 - 22). Everything else, including the subroutine *warpReduce* remains equal.

Listing 10: Kernel Optimized GPU Count Left

```

1  template <unsigned int blockSize>
2  extern __global__ void reduce2(
3  float * g_idata ,
4  unsigned int * g_begins ,
5  unsigned int * g_ends ,
6  float * g_cuts ,
7  unsigned int * g_odata) {
8      __shared__ unsigned int s_data[blockSize];
9
10     unsigned int tid = threadIdx.x;
11     const unsigned int begin = g_begins[blockIdx.x];
12     const unsigned int end = g_ends[blockIdx.x];
13     const float cut = g_cuts[blockIdx.x];
14
15     unsigned int i = begin + tid;
16     s_data[tid] = 0;
17
18     // unaligned coalesced g memory access
19     while (i < end) {
20         s_data[tid] += (g_idata[i] <= cut);
21         i += blockSize;
22     }
23     __syncthreads();
24
25     if (blockSize >= 512) {
26         if (tid < 256) {
27             s_data[tid] += s_data[tid + 256];
28         }
29         __syncthreads();
30     }
31     if (blockSize >= 256) {
32         if (tid < 128) {
33             s_data[tid] += s_data[tid + 128];
34         } __syncthreads();
35     }
36     if (blockSize >= 128) {
37         if (tid < 64) {
38             s_data[tid] += s_data[tid + 64];
39         } __syncthreads();
40     }
41     if (tid < 32) {
42         warpReduce<blockSize>(s_data, tid);
43     }
44     if (tid == 0) {
45         g_odata[blockIdx.x] = s_data[0];
46     }
47 }

```

6.6 Improved GPU Accelerated Count Left with Warp Level Primitives

The above described kernel code can also be rewritten where warp level primitives are leveraged to the full extent. Lines 1-14 are identical as in the previous kernel as

proposed in listing 10. On lines 16 and 17 we fetch the lane *id* and the *warpId*. The *lane* defines the index of a thread inside each warp. The operation $tid \& (32 - 1)$ is equivalent to $tid \% 32$ but significantly faster [NVIDIA, 2022d]. When incrementally assigning an index to each warp the index can be fetched by dividing the *tid* by the number of threads in a warp, which is then implicitly floored when converted into an integer. The operation $tid >> 5$ is equivalent to $tid / 32$ [NVIDIA, 2022d].

Instead of iterating over elements and storing them in global memory, the data is stored in local memory as seen on lines 19 to 24. Furthermore, at this stage there is no need to sync all threads within the block but only the threads within a warp, therefore the command `_syncthreads()` can be replaced with `_syncwarp()`. On lines 26 - 27 a warp level reduction is performed, which is equivalent to the one we have explained with the previous kernel. Now each thread with a *lane* equivalent to 0 has stored the sum of 31 next consecutive threads in local memory. In the case of 256 threads this means that there are 8 elements left to be summed up, however this step cannot be performed on a warp level as all values are stored in distinct warps. Therefore *val* is read from local memory and written to shared memory at an index of *warpId* and a `_syncthreads()` call is made (line 29-32). In this case it is necessary to synchronize all threads, as the storing of elements happens across warps. The number of values stored in shared memory can be determined by dividing the number of threads in a block by the warp size. Thus all threads for which a value in shared memory was written in the previous step, load the said value back to local memory, the others store 0 in their local memory (line 34-35). At this point a final warp level reduction can be performed (lines 34-41). A single element on thread 0 stores the entire sum across all input data elements and it can be written to global memory and later fetched from the CPU.

Listing 11: Kernel Optimized GPU Count Left

```

1 template <unsigned int blockSize>
2 extern __global__ void reduce3(
3 float * g_idata ,
4 unsigned int * g_begins ,
5 unsigned int * g_ends ,
6 float * g_cuts ,
7 unsigned int * g_odata) {
8     __shared__ unsigned int s_data[32];
9     unsigned int tid = threadIdx.x;
10    const unsigned int begin = g_begins[blockIdx.x];
11    const unsigned int end = g_ends[blockIdx.x];
12    const float cut = g_cuts[blockIdx.x];
13    unsigned int i = begin + tid;
14    s_data[tid] = 0;
15
16    const unsigned int lane = tid & (32 - 1);
17    const unsigned int warpId = tid >> 5;
18
19    unsigned int val = 0;
20    while (i < end) {
21        val += (g_idata[i] <= cut);
22        i += blockSize;
23    }
24    __syncwarp();
25
26    for (int offset = 16; offset > 0; offset >>= 1)
27        val += __shfl_down_sync(FULLMASK, val, offset);
28
29    if (lane == 0) {
30        s_data[warpId] = val;
31    }
32    __syncthreads();
33
34    if (warpId == 0) {
35        val = (tid < blockDim.x / warpSize) ? s_data[lane] : 0;
36
37        for (int offset = 16; offset > 0; offset >>= 1)
38            val += __shfl_down_sync(FULLMASK, val, offset);
39
40        if (tid == 0) {
41            g_odata[blockIdx.x] = val;
42        }
43    }
44 }

```

Against the expectation that reduction3 from listing ?? might outperform reduction2 from listing ??: Test results from Piza Daint of the two kernels showed very similar runtimes for both kernels. For a test set of 2^{27} uniformly random distributed particles, both kernels took 1.9 milliseconds to execute. From this figure an effective performance of 560 GB/s can be evaluated using the formula $2^{27} \times 4 \times 10^6 / 1.9$ [Harris, 2020b]. In this case 4 is the number of bytes which are loaded per particle po-

sition⁶ and the factor 10^6 is used to convert milliseconds to seconds. 560 GB/s comes close to the effective hardware limit of 732 GB/s and *kernel2* was already performing near optimally, thus the improvements made in *kernel3* have a negligible impact.

6.7 GPU Accelerated Partitioning

A GPU accelerated Count Left service is implemented successfully. Since the particles are still partitioned on the CPU, the data needs to be copied to GPU each time the SPTDS is expanded. As the Make Axis service is used to compress the axes into a single array, only N elements need to be copied, but this needs to be repeated $\lceil \log_2 d \rceil$ times. If the particles are partitioned on the GPU, it needs to be copied only once, but in this case all three coordinates have to be copied, thus $N \times 3$. Nevertheless this corresponds to a fraction of the total data transfer: $\frac{3}{\lceil \log_2 d \rceil \times}$. In figure 24 it is observable that a GPU partition becomes feasible, at least in terms of data transfer, after d exceeds 16.

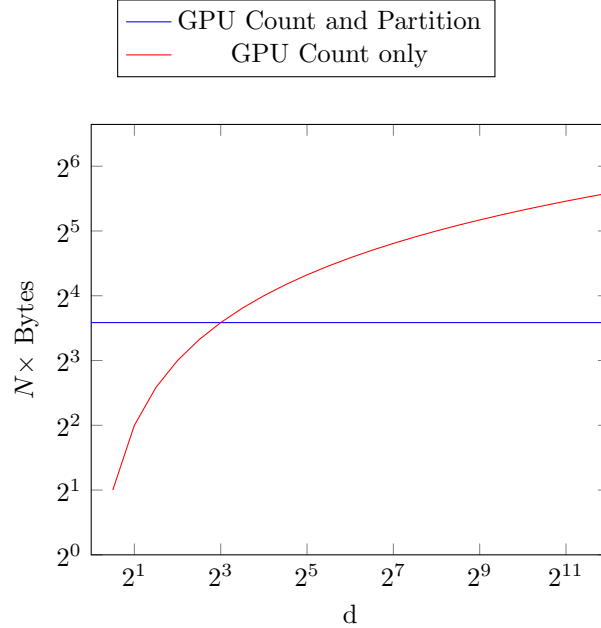


Figure 24: Data transferred with and without GPU partitioning

A GPU partitioning algorithm was first proposed by Cederman [Cederman and Tsigas, 2010] and later refined by Manca [Manca et al., 2016] where atomic operations are leveraged to improve the runtime of the algorithm. This version is built around the concepts introduced by Manca.

Early results have shown partitioning the data on the GPU is not really faster than the CPU version where the CPU version is not even implemented in a optimal

⁶We only iterate over a single axis

manner and finally the number of particles which can be stored on the GPU memory is reduced by $\frac{5}{3}$. This is caused by required additional temporary arrays for memory swaps on the GPU. Furthermore there was not enough time to fully optimize and track down all issues with the code, therefore not full kernel code is provided, instead the underlying concepts are introduced for subsequent work.

6.7.1 Schedule

In the adapted schedule (figure 25) the Copy Particles To GPU service can be moved out of the main loop. Other than that the schedule remains the same.

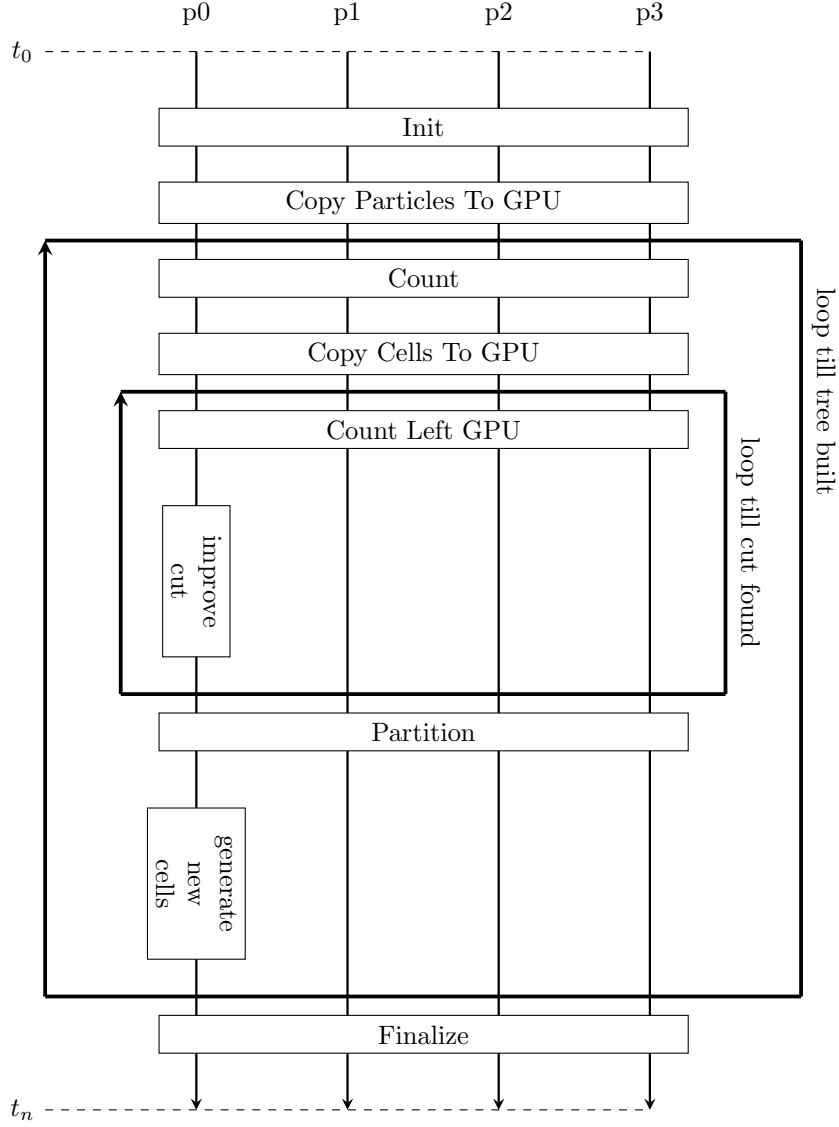


Figure 25: Master later schedule of a GPU accelerated ORB implementation with GPU partitioning

6.7.2 Partitioning a Single Cell

In a first step a block wise partition is performed, meaning there is no synchronization between blocks at this stage. Each thread loads the particle coordinate of the matching axis into the shared memory register S . Then two copies of the data stored in the shared memory are made where the first copy A compares each value with the cut

and stores 0 if the value is smaller and 1 otherwise. Consequently the second copy B stores the inverted result. Now a parallel prefix sum⁷ is computed on both copies, where an algorithm as described in GPU Gems [NVIDIA, 2022a] can be used. In B each element is then incremented by an offset equal to the total number of elements smaller than the cut. In theory the block wise partition could already be applied by storing each element with index i larger than the cut at A_i and larger elements at B_i . However the goal is to partition all elements across multiple blocks. To do so a global variable called leq and g is introduced, where the total number of elements smaller equal and bigger than the cut are stored. Each time a block has computed both scans, the number of elements designated to the block smaller and larger than the cut are known. Thus a single thread of each block performs an atomic add on leq and g accordingly. *AtomicAdd()* returns the old value of the variable. The return value can be caught and stored in two designated shared memory registers s_{leq} and s_g . At this stage every thread in a block agrees on a value for s_{leq} and s_g the values can be treated as starting indices to store the block wise partitioned results in the output array. In this case however an element i , which is larger than the cut is stored location $A_i + s_{leq}$ and if larger at $B_i = s_g$.

In this case the input array are coordinate values of particles and the output array is a temporary array of size N . The in an output array cannot be the same as the partition cannot be performed in place. Doing so would result in elements that are overwritten before they are read by the thread assigned to it. This cannot be avoided with the proposed technique as the relationship between the range of input data treated by a kernel and the output is not deterministic. Its dictated by the order, by which the kernel reaches the *atomicAdd()* statement.

As of now only a single axis is partitioned, but the same needs to be done for the other two axes. One approach is to partition all the other arrays in the same kernel, as the permutations can simply be reapplied. However in this case three additional temporary arrays are needed, one for each axis. Further Reducing the number of particles which can be stored on the GPU by $\frac{1}{2}$. A more memory efficient, but slightly slower variant is to allocate another array on the GPU where the actual permutation indices are stored when the partition kernel is executed. Therefore together with the temporary array N is only reduced by $\frac{3}{5}$. After the partition kernel has finished, the coordinates can be copied back from the temporary array to the original location in a device to device copy. Then the second axis is partitioned using the precomputed permutations and the results are again copied back to the proper location. Finally after the step is repeated for the last axis, the partition is finished.

6.7.3 Partitioning Multiple Cells

The same method as introduced in section 6.5 is applied to ensure only a single kernel has to be launched for all cells. Each block is assigned a cell, where all relevant information about the cell is stored in different arrays in global memory. Each block then reads the designated data by accessing the arrays at index *blockIdx.x*. The variables leq and g are replaced by arrays of the same length as there are cells to be partitioned⁸. This ensures all the blocks assigned to a single cell perform the atomic operations on the same variable. For a block to know which element in the leq and g arrays it needs to access, a sequential cell index array is passed, where all blocks working on the same cell receive the same cell index.

⁷Commonly referred to as a scan

⁸which is not equivalent to the total number of blocks

7 Performance Analysis of ORB

All measurements were performed on Piz Daint where the performance of the system can be read from table 1.

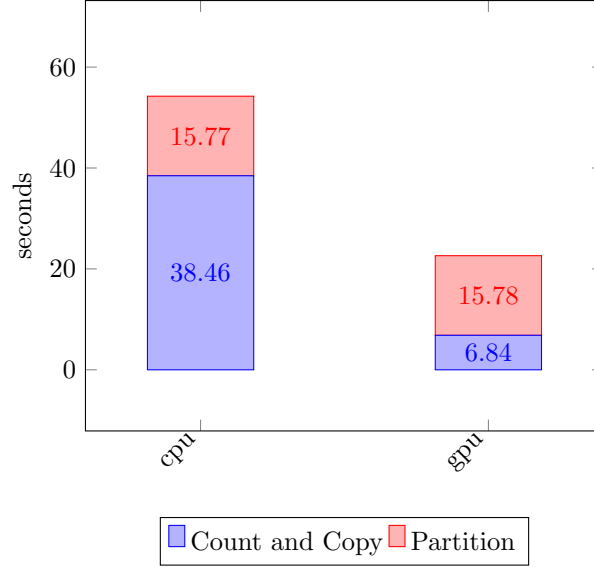


Figure 26: Runtimes for $N = 2^{30}$ and $d = 4096$ on a single Piz Daint node

In figure 26 the final measurements are plotted for $N = 2^{30}$ and $d = 4096$ on a single Piz Daint node. The reason why the measurements were not performed on multiple nodes, is that the paper does not propose any changes to the internode communication strategies, therefore these results are not of importance. A uniform random distribution was used instead of the target data, because at time of estimating the runtime there was no dataset available with the proper distribution. The results should however be comparable as solely the exact number of iterations taken to find a cut for each cell varies on different datasets.

A large portion of the total ORB runtime is attributed to the partition part. If the time lost during the partitioning is subtracted, the GPU version is about 5.6 times faster than the CPU implementation, which is slightly worse than the expected speedup of 6.2. The difference can be explained by the GPU reduction kernel which is slightly worse than the memory bandwidth would allow as described in section 6.6. The theoretical model on the other assumes an implementation which performs just as well as the memory bandwidth would allow.

When comparing the absolute runtimes there is a stark difference visible: 67.77s were predicted for the CPU version where 38.46s was the actual timing. Similarly the GPU was expected to execute ORB⁹ in 10.79s, however it only took 6.84s. The theoretical model considers an upper bound for the total number of iterations which the bisection method takes to converge. In reality this number is closer to 20 than to p

⁹Excluding the time required for the partition

which is 32 in the input data. When considering this, the measurements are also in an absolute term relatively close to the theoretical model. $67.47s \times \frac{20}{32} = 42.16s$ against the actual $38.46s$ and $10.79s \times \frac{20}{32} = 6.74s$ compared to the measured $6.84s$. Concluding the measured runtime are surprisingly similar to the expected ones, however they do not reach the theoretical maximum. The reduction kernel would need to be further improved to fully leverage the GPU bandwidth, however as of writing this paper no better implementation variant is known.

7.1 Scaling

The CPU and GPU versions were also tested for different problem sizes N . The results were then compiled into speedups, enabling us to see how the speedup behaves for different problem sizes.

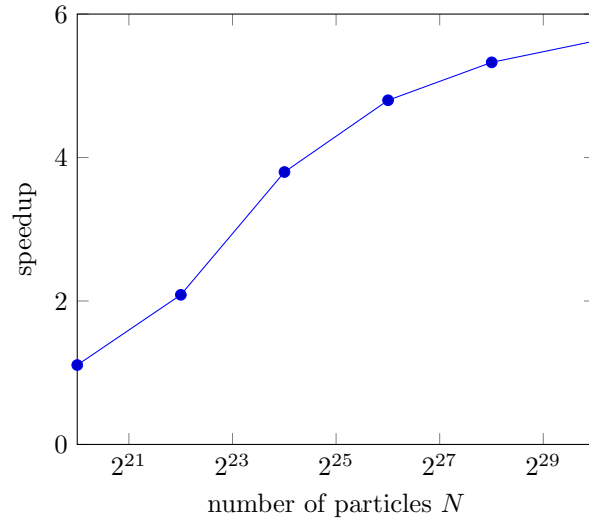


Figure 27: Speedup of GPU over CPU version for different problem sizes and $d = 4096$

In the resulting figure 27 it can be observed how the speedup is only minor for *small* problem sizes. This is most certainly caused by an overhead introduced by the GPU implementation over its CPU counterpart, as kernel invocations and memory copying operations are not free of cost. As soon as the problem size approaches a more realistic size, meaning closer to the number of particles which are to be expected on a single node, the speedup appears to converge to a factor of 6. This is a positive outcome, as it would be undesirable to have a degrading performance for larger problem sizes.

8 Conclusion

This research aimed at improving the runtime of the ORB algorithm, which constructs a space partitioning tree data structure, using the GPU programming language CUDA. The SPTDS is used for the FMM method as well as load balancing.

In order to have an idea how large the speedup of an optimal GPU implementation against an optimal CPU version is, I have calculated runtime for both versions. In theory a maximal speedup of 6.2 was expected.

With extensive research about the CUDA programming language I have written three GPU reduction kernels, where kernel2 significantly outperformed the first one. The kernel essentially solves a segmented map reduce problem. NVIDIA's CUB [NVIDIA, 2022c] exposes a similar function, however in the case of ORB, the map condition varies across segments which is not supported in CUB. I have introduced another version, kernel3 which makes extensive use of warp level primitives and can reduce access to shared memory by a factor of 32, however it did not outperform kernel2 and measurements have shown that the throughput of kernel2 and kernel3 is close to GPU memory bandwidth.

With the gained knowledge I have tried to accelerate the partitioning using CUDA. However due to time constraints I was not able to finish writing the required kernels. Thus I described the theoretical concepts, which can build a solid foundation for further research. However it is not clear whether the GPU partitioning could further improve the ORB runtime when implemented. Finally I have made time measurements on a single node of the Piz Daint supercomputer. When isolating the time it takes to compute a cut, without considering the time used for partitioning, the GPU accelerated version I have implemented runs about 5.6 times faster than a fully optimized and parallelized CPU version. Which comes close to the theoretical limit and the initial goal of improving the runtime of ORB is successful.

The section of ORB which could still be improved by a lot, is the choice of the bisection method. As of now many iterations are needed to find a ideal cut. In theory more sophisticated root finding methods or binning algorithms[Voglsam, 2013] could be employed, further improving the runtime of ORB.

References

- [AMD, 2022] AMD (2022). AMD EPYC™ 7742 | AMD. [Online; accessed 27. Apr. 2022].
- [Cederman and Tsigas, 2010] Cederman, D. and Tsigas, P. (2010). GPU-Quicksort: A practical Quicksort algorithm for graphics processors. *ACM J. Exp. Algorithmics*, 14:4–1.
- [Cormen et al., 2022] Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2022). *Introduction to algorithms*. MIT press.
- [CSCS, 2022] CSCS (2022). Piz Daint & Piz Dora. [Online; accessed 27. Apr. 2022].
- [Harris, 2020a] Harris, M. (2020a). How to Access Global Memory Efficiently in CUDA C/C++ Kernels | NVIDIA Technical Blog. [Online; accessed 11. Aug. 2022].
- [Harris, 2020b] Harris, M. (2020b). How to Implement Performance Metrics in CUDA C/C++ | NVIDIA Technical Blog. [Online; accessed 12. Aug. 2022].
- [Harris, 2021a] Harris, M. (2021a). How to Optimize Data Transfers in CUDA C/C++ | NVIDIA Technical Blog. [Online; accessed 12. Aug. 2022].
- [Harris, 2021b] Harris, M. (2021b). Using Shared Memory in CUDA C/C++ | NVIDIA Technical Blog. [Online; accessed 11. Aug. 2022].
- [Harris, 2022] Harris, M. (2022). Optimizing parallel reduction in cuda. [Online; accessed 1. Aug. 2022].
- [Intel, 2022] Intel (2022). Intel® Xeon® Processor E5-2690 . [Online; accessed 27. Apr. 2022].
- [Manca et al., 2016] Manca, E., Manconi, A., Orro, A., Armano, G., and Milanese, L. (2016). CUDA-quicksort: an improved GPU-based implementation of quicksort. *Concurrency Computat.: Pract. Exper.*, 28(1):21–43.
- [NVIDIA, 2022a] NVIDIA (2022a). Chapter 39. Parallel Prefix Sum with CUDA. [Online; accessed 3. Aug. 2022].
- [NVIDIA, 2022b] NVIDIA (2022b). CUB. [Online; accessed 28. Jul. 2022].
- [NVIDIA, 2022c] NVIDIA (2022c). CUB: DeviceSegmentedReduce Struct Reference. [Online; accessed 7. Aug. 2022].
- [NVIDIA, 2022d] NVIDIA (2022d). CUDA C++ Best Practices Guide. [Online; accessed 11. Aug. 2022].
- [NVIDIA, 2022e] NVIDIA (2022e). CUDA C++ Programming Guide. [Online; accessed 2. Aug. 2022].
- [NVIDIA, 2022f] NVIDIA (2022f). NVIDIA Tesla P100: The Most Advanced Data Center Accelerator. [Online; accessed 27. Apr. 2022].
- [NVIDIA, 2022g] NVIDIA (2022g). Thrust. [Online; accessed 28. Jul. 2022].
- [ORNL, 2022] ORNL (2022). ORNL Launches Summit Supercomputer | ORNL. [Online; accessed 27. Apr. 2022].
- [Rennich, 2021] Rennich, S. (2021). CUDA C/C++ Streams and Concurrency. [Online; accessed 12. Aug. 2022].
- [Stadel, 2001] Stadel, J. G. (2001). *Cosmological N-body simulations and their analysis*. PhD thesis.

- [TOP500, 2022] TOP500 (2022). TOP500. [Online; accessed 1. Aug. 2022].
- [UZH, 2022] UZH (2022). ICS Software. [Online; accessed 1. Aug. 2022].
- [Veldhuizen, 2022] Veldhuizen, T. (2022). Blitz++ User’s Guide . [Online; accessed 11. May 2022].
- [Voglsam, 2013] Voglsam, G. (2013). Real-time ray tracing on the gpu - ray tracing using cuda and kd-trees. Master’s thesis, Institute of Computer Graphics and Algorithms, Vienna University of Technology.
- [Yuan Lin, 2021] Yuan Lin, V. G. (2021). Using CUDA Warp-Level Primitives | NVIDIA Technical Blog. [Online; accessed 11. Aug. 2022].

List of Figures

1	Uniform random distribution of 3D coordinates in cube domain projected onto a 2d plane	9
2	Example 2D particle distribution with $N = 7$	13
3	Volumes of SPTDS leaf cells generated from ORB at a depth of 0 . . .	13
4	SPTDS generated from ORB at a recursion depth of 0	13
5	Volumes of SPTDS leaf cells generated from ORB at a depth of 1 . . .	14
6	SPTDS generated from ORB at a recursion depth of 1	14
7	Volumes of SPTDS leaf cells generated from ORB at a depth of 2 . . .	15
8	SPTDS generated from ORB at a recursion depth of 2	15
9	Computing model	21
10	Roofline Model for the isolated Count Left code on Piz Daint	25
11	Roofline Model for the isolated Count Left code on Piz Daint using the GPU	25
12	Roofline Model for the isolated Count Left code on Piz Daint	26
13	Roofline Model for the isolated Count Left code on Summit using the GPU	26
14	Roofline Model for the isolated Count Left code on Alps	27
15	Estimated runtimes of different variants	31
16	Tree stored as an array heap	33
17	SPTDS constructed from example particle distribution tree with 3 leaf cells ($d = 3$)	34
18	Example particles distributed randomly across 3 nodes	39
19	Master later schedule of a CPU only ORB implementation	40
20	Fundamental CUDA concepts	41
21	Master later schedule of a GPU accelerated ORB implementation . . .	46
22	Kernel calls with GPU Count Left with two different versions	50
23	Master later schedule of an improved GPU accelerated ORB implementation	51
24	Data transferred with and without GPU partitioning	56
25	Master later schedule of a GPU accelerated ORB implementation with GPU partitioning	58
26	Runtimes for $N = 2^{30}$ and $d = 4096$ on a single Piz Daint node	60
27	Speedup of GPU over CPU version for different problem sizes and $d = 4096$	61

List of Tables

1	Datapoints of Supercomputers	22
2	Variables to estimate gflops	23
3	Variables used for runtime estimates	28